UNIFYING CAUSAL REPRESENTATION LEARNING WITH THE INVARIANCE PRINCIPLE

Anonymous authors Paper under double-blind review

ABSTRACT

Causal representation learning aims at recovering latent causal variables from high-dimensional observations to solve causal downstream tasks, such as predicting the effect of new interventions or more robust classification. A plethora of methods have been developed, each tackling carefully crafted problem settings that lead to different types of identifiability. The folklore is that these different settings are important, as they are often linked to different rungs of Pearl's causal hierarchy, although not all neatly fit. Our main contribution is to show that many existing causal representation learning approaches methodologically align the representation to known data symmetries. Identification of the variables is guided by equivalence classes across different "data pockets" that are not necessarily causal. This result suggests important implications, allowing us to unify many existing approaches in a single method that can mix and match different assumptions, including non-causal ones, based on the invariances relevant to our application. It also significantly benefits applicability, which we demonstrate by improving treatment effect estimation on real-world high-dimensional ecological data. Overall, this paper clarifies the role of causality assumptions in the discovery of causal variables and shifts the focus to preserving data symmetries.

027 1 INTRODUCTION

004

005

010

011

012

013

014

015

016

017

018

019

021

025

026

Causal representation learning (Schölkopf et al., 2021) posits that many real-world high-dimensional 029 perceptual data can be described through a simplified latent structure specified by a few interpretable 030 low-dimensional causally-related variables. Discovering hidden causal structures from data has been 031 a long-standing goal across many scientific disciplines, spanning neuroscience (Vigário et al., 1997; Brown et al., 2001), communication theory (Ristaniemi, 1999; Donoho, 2006), economics (Angrist 033 & Pischke, 2009) and social science (Antonakis & Lalive, 2011). From the machine learning per-034 spective, algorithms and models integrated with causal structure are often proven to be more robust at distribution shift (Ahuja et al., 2022a; Bareinboim & Pearl, 2016; Rojas-Carulla et al., 2018), providing better out-of-distribution generalization results and reliable agent planning (Fumero et al., 037 2024; Seitzer et al., 2021; Urpí et al., 2024). The general goal of causal representation learning ap-038 proaches is formulated as to provably identify ground-truth latent causal variables and their causal relations (up to certain ambiguities). Many existing approaches in causal representation learning carefully formulate their problem settings to guarantee identifiability and justify the assumptions 040 within the framework of Pearl's causal hierarchy, such as "observational, interventional, or coun-041 terfactual CRL" (von Kügelgen et al., 2024; Ahuja et al., 2023; Brehmer et al., 2022; Buchholz 042 et al., 2024; Zhang et al., 2024a; Varici et al., 2024a). However, some causal representation learning 043 works may not perfectly fit within this causal language framework; for instance, the problem set-044 ting of temporal CRL works (Lachapelle et al., 2022; Lippe et al., 2022a;; 2023) does not always 045 align straightforwardly with existing categories. They often assume that an individual trajectory is 046 "intervened" upon, but this is not an intervention in the traditional sense, as noise variables are not 047 resampled. It is also not a counterfactual, as the value of non-intervened variables can change due 048 to default dynamics. Similarly, domain generalization (Sagawa et al., 2019; Krueger et al., 2021; Ahuja et al., 2022a) and certain multi-task learning approaches (Lachapelle et al., 2023; Fumero et al., 2024) are sometimes framed as informally related to causal representation learning. How-051 ever, the precise relation to causality is not always clearly articulated. This has resulted in a variety of methods and findings, some of which rely on assumptions that might be too narrowly tailored 052 for practical, real-world applications. For example, Cadei et al. (2024) collected a data set for estimating treatment effects from high-dimensional observations in real-world ecology experiments.

054 Despite the clear causal focus of the benchmark, they note that even having access to multiple views 055 and being able to perform some interventions, neither existing multi-view nor interventional causal 056 representation learning methods are directly applicable due to mismatching assumptions. 057

This paper contributes a unified rephrasing of many existing nonparametric CRL works through the 058 lens of invariance (Peters et al., 2014; Heinze-Deml et al., 2018; Arjovsky et al., 2020). We observe that many existing causal representation approaches share methodological similarities, particularly 060 in aligning the representation with known data symmetries, while differing primarily in how the 061 invariance principle is invoked. This invariance principle is usually formulated *implicitly* in the 062 assumed data-generating process. Instead, we make this explicit and show that latent causal variable 063 identification broadly originates from multiple data pockets with certain underlying equivalence 064 relations. These are not necessarily causal and (with minor exceptions) have to be known apriori. 065 Unifying causal representation learning approaches using the invariance principle brings several 066 potential benefits: First, it helps clarify the alignment between seemingly different categories 067 of CRL methods, contributing to a more coherent and accessible framework for understanding 068 causal representation learning. This perspective may also allow for the integration of multiple invariance relations in latent variable identification, which could improve the flexibility of these 069 methods in certain practical settings. Additionally, our theory underscores a gap in the necessary causal assumptions for graph learning, which is essential for generalizing to unseen interventions 071 and helps distinguish it from the problem of identifying causal variables. These invariances can 072 be expressed in causal terms, such as interventions, but do not always need to be. Last but not 073 least, this formulation of invariance relation links causal representation learning to many existing 074 representation learning areas outside of causality, including invariant training (Arjovsky et al., 2020; 075 Ahuja et al., 2022a), domain adaptation (Sagawa et al., 2019; Krueger et al., 2021), and geometric 076 deep learning (Cohen & Welling, 2016; Bronstein et al., 2017; 2021). 077

We highlight our contributions as follows: 078

079 • We propose a unified rephrasing for existing nonparametric causal representation learning approaches leveraging the invariance principles and proving latent variable identifiability in this 081 general setting (§ 3). We show that 30 existing identification results can be seen as special 082 cases directly implied by our framework (Tab. 4). This approach also enables us to derive new 083 results, including latent variable identifiability from one imperfect intervention per node in the 084 nonparametric setting (Cor. D.1). 085

- In addition to employing different methods, many works in the CRL literature use varying definitions of "identifiability." We formalize these definitions at different levels of granularity, highlight their connections, and demonstrate how various definitions can be addressed within our framework by leveraging different invariance principles (App. C.1).
- Upon the identifiability of the latent variables, we discuss the necessary causal assumptions for graph identification and the possibility of partial graph identification using the language of causal 092 consistency. With this, we draw a distinction between the causal assumptions necessary for graph discovery and those that may not be required for variable discovery (App. C.2). 093
 - Our framework is broadly applicable across a range of settings. We observe improved results on real-world experimental ecology data using the causal inference benchmark from highdimensional observations provided by Cadei et al. (2024) (§ 5.1). Additionally, we present a synthetic ablation to demonstrate that existing methods, which assume access to interventions, actually only require a form of distributional invariance to identify variables. This invariance does not necessarily need to correspond to a valid causal intervention (\S 5.2).
 - 2 PROBLEM SETTING

087

088

089

091

094

095

096

097

098

099

100 101

102

103 **Notation.** [N] is used as a shorthand for $\{1, \ldots, N\}$. We use bold lower-case z for random vectors 104 and normal lower-case z for their realizations. A vector z can be indexed either by a single index $i \in$ 105 $[\dim(\mathbf{z})]$ via \mathbf{z}_i or a index subset $A \subseteq [\dim(\mathbf{z})]$ with $\mathbf{z}_A := \{\mathbf{z}_i : i \in A\}$. $P_{\mathbf{z}}$ denotes the probability distribution of the random vector z and $p_z(z)$ denotes the associated probability density function. 106 By default, a "measurable" function is *measurable* w.r.t. the Borel sigma algebras and defined w.r.t. 107 the Lebesgue measure. A more comprehensive summary of notations is provided in App. A.

144

145

146

147

156

159

160

161

Category	Example	Invariance
Multiview CRL	$ \begin{array}{c} \mathbf{z}_1^1 \mathbf{z}_2^1 \\ \mathbf{v} \end{array} \mathbf{z}_3^1 \mathbf{z}_4^1 \\ \mathbf{z}_1^2 \mathbf{z}_2^2 \mathbf{z}_3^2 \mathbf{z}_4^2 \end{array} $	Sample level invariance
		$\mathbf{z}_A^1 = \mathbf{z}_A^2$
Multi-env. CRL (two	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Invariance on the
interventions per	$\mathbf{v}_{\mathbf{v}^1}$	interventional target
node)		$\mathcal{I}_{\mathbf{z}_A^1} = \mathcal{I}_{\mathbf{z}_A^2} = \{1\}$
	$\begin{bmatrix} 1 \\ r^1 \end{bmatrix} \begin{bmatrix} 1 \\ r^1 \end{bmatrix} \begin{bmatrix} 1 \\ r^2 \end{bmatrix} \begin{bmatrix} r^2 $	
Multi-env. CRL (one	\mathbf{z}_1 \mathbf{z}_2 \mathbf{z}_3 \mathbf{z}_4 \mathbf{z}_1 \mathbf{z}_2 \mathbf{z}_3 \mathbf{z}_4	Marginal invariance
intervention per	(\mathbf{x}^1)	$p_{\mathbf{z}_A^1} = p_{\mathbf{z}_A^2}$
node)	<u> </u>	
Multi-env CRL (one	(\mathbf{z}_1^1) (\mathbf{z}_2^1) (\mathbf{z}_3^1) (\mathbf{z}_4^1) (\mathbf{z}_1^2) (\mathbf{z}_2^2) (\mathbf{z}_3^2) (\mathbf{z}_4^2)	Score invariance
intervention per		$S_{\pi^1} = S_{\pi^2}$
node)	$\begin{pmatrix} \mathbf{x}^1 \end{pmatrix}$ $\begin{pmatrix} \mathbf{x}^2 \end{pmatrix}$	\mathbf{z}_A \mathbf{z}_A
Temporal CRL	$\begin{pmatrix} \tilde{\mathbf{z}}_1^t \\ 1 \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{z}}_2^t \\ 2 \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{z}}_3^t \\ 1 \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{z}}_4^t \\ 2 \end{pmatrix} \begin{pmatrix} \mathbf{z}_1^t \\ 1 \end{pmatrix} \begin{pmatrix} \mathbf{z}_2^t \\ 2 \end{pmatrix} \begin{pmatrix} \mathbf{z}_3^t \\ \mathbf{z}_4^t \end{pmatrix}$	Transition invariance
-		$p_{\mathbf{z}_A \mathbf{z}_{t-1}} = p_{\tilde{\mathbf{z}}_A \mathbf{z}_{t-1}}$
	$\begin{pmatrix} \mathbf{y}^{\mathbf{r}} \\ \mathbf{x} \end{pmatrix} \begin{pmatrix} \mathbf{y}^{\mathbf{r}} \\ \mathbf{y}^{\mathbf{r}} \end{pmatrix}$	
Multi-task CRL	T_1 (\mathbf{z}_1) (\mathbf{z}_2) (\mathbf{z}_3) (\mathbf{z}_4) T_2	Overlapping task support
		$\mathbf{z}_A^{{\scriptscriptstyle I}_1} = \mathbf{z}_A^{{\scriptscriptstyle I}_2}$
	$\begin{pmatrix} \mathbf{y}^1 \end{pmatrix}$ $\begin{pmatrix} \mathbf{y}^2 \end{pmatrix}$	
Domain	\mathbf{z}_1^1 \mathbf{z}_2^1 \mathbf{z}_2^1 \mathbf{z}_2^1 \mathbf{z}_2^1 \mathbf{z}_2^2 \mathbf{z}_2^2 \mathbf{z}_2^2	Risk invariance on optimal
generalization		weights $\mathcal{R}_1^*(\mathbf{w}^*\mathbf{z}_A^1, \mathbf{y}^1) =$

Table 1: Examples of different CRL categories and their corresponding invariance. The invariant partition A is highlighted with a smoke blue box.

This section defines our problem setting using standard CRL concepts and assumptions (Formal definitions are deferred to App. B). While prior works in CRL typically categorize their settings using established causal language (such as "counterfactual," "interventional," or "observational"), our approach introduces a more general invariance principle that aims to unify diverse problem settings. We introduce the following concepts as mathematical tools to describe our data generating process.

Definition 2.1 (Invariance property). Let $A \subseteq [N]$ be an index subset of the Euclidean space \mathbb{R}^N and let \sim_{ι} be an equivalence relationship on $\mathbb{R}^{|A|}$, with A of known dimension. Let $\mathcal{M} := \mathbb{R}^{|A|} / \sim_{\iota}$ be the quotient of $\mathbb{R}^{|A|}$ under this equivalence relationship; \mathcal{M} is a topological space equipped with the quotient topology. Let $\iota : \mathbb{R}^{|A|} \to \mathcal{M}$ be the projection onto the quotient induced by the equivalence relationship \sim_{ι} . This projection ι is termed the *invariance property* of this equivalence relation. Two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^{|A|}$ are invariant under ι if and only if they belong to the same \sim_{ι} equivalence class, i.e.:

$$\iota(\mathbf{a}) = \iota(\mathbf{b}) \Leftrightarrow \mathbf{a} \sim_{\iota} \mathbf{b}.$$

Extending this definition to the whole latent space \mathbb{R}^N , a pair of latents $\mathbf{z}, \tilde{\mathbf{z}} \in \mathbb{R}^N$ are *non-trivially invariant on a subset* $A \subseteq [N]$ *under the property* ι only if

(i) the invariance property ι holds on the indices $A \subseteq [N]$ in the sense that $\iota(\mathbf{z}_A) = \iota(\tilde{\mathbf{z}}_A)$;

(ii) for any smooth functions $h_1, h_2 : \mathbb{R}^N \to \mathbb{R}^{|A|}$, the invariance property between $\mathbf{z}, \tilde{\mathbf{z}}$ breaks under the h_1, h_2 transformations if h_1 or h_2 directly depends on some other component \mathbf{z}_q with $q \in [N] \setminus A$. Taking h_1 and \mathbf{z} as an example, we have:

$$\exists q \in [N] \setminus A, \mathbf{z}^* \in \mathbb{R}^N, \quad s.t. \ \frac{\partial h_1}{\partial \mathbf{z}_q}(\mathbf{z}^*) \text{ exists and is non zero} \quad \Rightarrow \quad \iota(h_1(\mathbf{z})) \neq \iota(h_2(\tilde{\mathbf{z}}))$$

which means: given that the partial derivative of h_1 w.r.t. some latent variable $\mathbf{z}_q \in \mathbf{z}_{[N]\setminus A}$ exists and is non-zero at some point $\mathbf{z}^* \in \mathbb{R}^N$, $h_1(\mathbf{z})$, $h_2(\mathbf{z})$ violates the invariance principle in the sense that $\iota(h_1(\mathbf{z})) \neq \iota(h_2(\tilde{\mathbf{z}}))$.

Intuition: The invariance property ι maps the invariant latent subset \mathbf{z}_A to the space \mathcal{M} representing the identified factor of variations. For example, in the multi-view literature (von Kügelgen et al., 2021; Brehmer et al., 2022; Yao et al., 2023), it is the *identity map* because the pre-and post action views are sharing the *exact value* of the invariant latents; for the interventional and temporal CRL (Varici et al., 2023; von Kügelgen et al., 2024; Lachapelle et al., 2022; Lippe et al., 2022a), this invariance property holds on a *distributional* level, and the property manifold \mathcal{M} can play the role of parameter space for the parameteric latent distributions or the general distribution space for the nonparametric case; for the multi-task line of work (Lachapelle et al., 2023; Fumero et al., 2024), ι maps the task-related latents to the overlapping task support. Concrete examples of each special case, along with the explicit formulation of their invariance, are provided in Tab. 1.

180 **Remark**: Defn. 2.1 (ii) is essential for latent variable identification on the invariant partition A, 181 which is further justified in App. E.1 by showing a non-identifiable example violating (ii). Intu-182 itively, Defn. 2.1 (ii) presents sufficient discrepancy between the invariant and variant part in the 183 ground truth generating process, paralleling various key assumptions for identifiability in CRL that were termed differently but conceptually similar, such as sufficient variability (von Kügelgen et al., 2024; Lippe et al., 2022b), interventional regularity (Varici et al., 2023; 2024b) and interventional 185 discrepancy (Wendong et al., 2024; Varici et al., 2024a). On a high level, these assumptions guar-186 antee that the intervened mechanism sufficiently differs from the default causal mechanism to effec-187 tively distinguish the intervened and non-intervened latent variables, which serves the same purpose 188 as Defn. 2.1 (ii). We elaborate this link further in App. E.1. 189

190 We denote by $S_{\mathbf{z}} := {\mathbf{z}^1, \dots, \mathbf{z}^K}$ the set of latent random vectors with $\mathbf{z}^k \in \mathbb{R}^N$ and write its 191 joint distribution as $P_{S_{\mathbf{z}}}$. The joint distribution $P_{S_{\mathbf{z}}}$ has a probability density $p_{S_{\mathbf{z}}}(z^1, \dots, z^K)$. Each 192 individual random vector $\mathbf{z}^k \in S_{\mathbf{z}}$ follows the marginal density $p_{\mathbf{z}^k}$ with the non-degenerate support 193 $\mathcal{Z}^k \subseteq \mathbb{R}^N$, whose interior is a non-empty open set of \mathbb{R}^N .

Definition 2.2 (Observable of a set of latent random vectors). Consider a set of random vectors $S_{\mathbf{z}} := {\mathbf{z}^1, \dots, \mathbf{z}^K}$ with $\mathbf{z}^k \in \mathbb{R}^N$, the corresponding set of observables $S_{\mathbf{x}} := {\mathbf{x}^1, \dots, \mathbf{x}^K}$ is generated by $S_{\mathbf{x}} = F(S_{\mathbf{z}})$, where the map *F* defines a push-forward measure $F_{\#}(P_{S_{\mathbf{z}}})$ on the image of *F* as:

$$F_{\#}(P_{\mathcal{S}_{\mathbf{z}}})(\mathbf{x}^{1},\dots,\mathbf{x}^{K}) = P_{\mathcal{S}_{\mathbf{z}}}(f_{1}^{-1}(\mathbf{x}^{1}),\dots,f_{K}^{-1}(\mathbf{x}^{K}))$$
(2.1)

with the support $\mathcal{X} := \text{Im}(f) \subseteq \mathbb{R}^{K \times D}$. Note that F satisfies the diffeomorphism assumption (Asm. B.1) as each f_k is a diffeomorphism onto its image according to Asm. B.1.

Intuition. Defn. 2.2 formulates the generating process of the set of observables as a joint 202 pushforward of a set of latent random vectors, providing a formal definition of the non-iid. data 203 pockets employed in causal representation learning algorithms. It conveniently explains various 204 underlying data symmetries given inherently by individual problem settings. For example, in 205 the multiview scenario (von Kügelgen et al., 2021; Daunhawer et al., 2023; Yao et al., 2023), 206 we can observe the joint data distribution P_{S_x} because the data are "paired" (non-independent). 207 In the interventional CRL that relies on multi-environment data, the joint data distribution 208 can be factorized as a product of individual non-identical marginals $\{P_{\mathbf{x}^k}\}_{k \in [K]}$, originating from partially different latent distributions P_{z^k} that are modified by, e.g., interventions. In the 210 supervised setting, such as multi-task CRL, we have an extended data pocket augmented by the task labels that is formally defined as $S_{\bar{\mathbf{x}}} := \{\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_K\}$ with $\bar{\mathbf{x}}_k := (\mathbf{x}, \mathbf{y}^k)$. Note that the observable \mathbf{x} is shared across all tasks $k \in [K]$ whereas the tasks labels \mathbf{y}^k are specific to 211 212 individual tasks, thus introducing different joint data-label distributions $P_{\bar{\mathbf{x}}_{k}}$. 213

214

201

162

163 164

166

167

169 170

171

172

173

174

175

176

177

178

179

In the following, we denote by $\mathfrak{I} := \{\iota_i : \mathbb{R}^{|A_i|} \to \mathcal{M}_i\}$ a finite set of invariance properties with their respective invariant subsets $A_i \subseteq [N]$ and their equivalence relationships \sim_{ι_i} , each inducing

a projection onto its quotient and invariance property ι_i (Defn. 2.1). For a set of observables $\mathcal{S}_{\mathbf{x}} := {\mathbf{x}^1, \dots, \mathbf{x}^K} \in \mathcal{X}$ generated from the data generating process described in § 2, we assume:

Assumption 2.1. For each $\iota_i \in \Im$, there exists a *unique known* index subset $V_i \subseteq [K]$ with at least two elements (i.e., $|V_i| \ge 2$) s.t. $\mathbf{x}_{V_i} = F([\mathbf{z}]_{\sim_{\iota_i}})$ (which we term informally as "data pockets") forms the set of observables generated from an equivalence class $[\mathbf{z}]_{\sim_{\iota_i}} := \{\tilde{\mathbf{z}} \in \mathbb{R}^N : \mathbf{z}_{A_i} \sim_{\iota_i} \tilde{\mathbf{z}}_{A_i}\}$, as given by Defn. 2.2.

Remark: Intuitively, Asm. 2.1 ensures that for each invariance property $\iota_i \in \mathfrak{I}$, there are at least two observables generated from latents that share ι_i ; otherwise the invariance partition A_i becomes undefined and no identification results can be derived. While \mathfrak{I} does not need to be fully described with explicit forms, which observables should belong to the same equivalence class is known (denoted as $V_i \subseteq [K]$ for the invariance property $\iota_i \in \mathfrak{I}$). This is a standard assumption and is equivalent to knowing, e.g., two views are generated from partially overlapped latents (Yao et al., 2023).

Problem setting. Given a set of observables $S_x \in \mathcal{X}$ satisfying Asm. 2.1, we show that we can simultaneously identify multiple invariant latent blocks A_i under a set of weak assumptions. In the best case, if each individual latent component is represented as a single invariant block through individual invariance property $\iota_i \in \mathcal{I}$, we can learn a fully disentangled representation and further identify the latent causal graph by additional technical assumptions.

3 IDENTIFIABILITY THEORY VIA THE INVARIANCE PRINCIPLE

239 High-level overview. This section presents a general theory for latent variable identification that 240 brings together many identifiability results from existing CRL works, including multiview, inter-241 ventional, temporal, and multi-task CRL. Our theory of latent variable identifiability, based on the 242 invariance principle, consists of two key components: (1) ensuring the encoder's sufficiency, thereby 243 obtaining an adequate representation of the original input for the desired task; (2) guaranteeing the 244 learned representation to preserve known data symmetries as invariance properties. The sufficiency 245 is often enforced by minimizing the reconstruction loss (Locatello et al., 2020; Ahuja et al., 2022b; Lippe et al., 2022b;a; Lachapelle et al., 2022) in auto-encoder based architecture, maximizing the 246 log likelihood in normalizing flows or maximizing entropy (Zimmermann et al., 2021; von Kügel-247 gen et al., 2021; Daunhawer et al., 2023; Yao et al., 2023) in self-supervised approaches. The in-248 variance property in the learned representations is often enforced by minimizing some equivalence 249 relation-induced regularizer (von Kügelgen et al., 2021; Yao et al., 2023; Lippe et al., 2022b; Zhang 250 et al., 2024a) or by some iterative algorithm that provably ensures the invariance property on the 251 output (Squires et al., 2023; Varici et al., 2024b). As a result, all invariant blocks A_i , $i \in [n_{\mathcal{I}}]$ can be 252 identified up to a mixing within the blocks while being disentangled from the rest. This type of iden-253 tifiability is defined as *block-identifiability* (von Kügelgen et al., 2021) which we restate as follows: 254

Definition 3.1 (Block-identifiability (von Kügelgen et al., 2021)). A subset $\mathbf{z}_A := {\mathbf{z}_j}_{j \in A}$ with $A \subseteq [N]$ of the latent variables is block-identified by an encoder $g : \mathbb{R}^D \to \mathbb{R}^N$ on the invariant subset A if the learned representation $\hat{\mathbf{z}}_{\hat{A}} := [g(\mathbf{x})]_{\hat{A}}$ with $\hat{A} \subseteq [N], |A| = |\hat{A}|$ contains all and only information about the ground truth \mathbf{z}_A , i.e. $\hat{\mathbf{z}}_{\hat{A}} = h(\mathbf{z}_A)$ for some diffeomorphism $h : \mathbb{R}^{|A|} \to \mathbb{R}^{|A|}$.

Intuition: Note that the inferred representation $\hat{z}_{\hat{A}}$ can be a set of entangled latent variables rather than a single one. Block-identifiability can be considered as a coarse-grained definition of disentanglement (Locatello et al., 2020; Fumero et al., 2024; Lachapelle et al., 2023), which seeks to disentangle individual latent factors. In other words, disentanglement can be considered a special case of block-identifiability, with each latent constituting a single invariant block. Notably, in (Locatello et al., 2020), disentangled factors were identified in blocks, with fine-grained identifiability achieved by intersecting different blocks.

260

261

262

264

265

224

225

226

227

228

229 230 231

232

233

234

235 236 237

238

Definition 3.2 (Encoders). The encoders $G := \{g_k : \mathcal{X}^k \to \mathcal{Z}^k\}_{k \in [K]}$ consist of smooth functions mapping from the observational support \mathcal{X}^k to the corresponding latent support \mathcal{Z}^k (§ 2).

275

276 277

278

279

280 281

282

283 284

285

286 287 288

289

290

295

313

Intuition: For the purpose of generality, we design the encoder g_k to be specific to individual observable $\mathbf{x}^k \in S_{\mathbf{x}}$. However, multiple g_k can share parameters if they work on the same modality. Ideally, we would like the encoders to preserve as much invariance (from \mathfrak{I}) as possible. Thus, a clear separation between different encoding blocks is needed. To this end, we introduce selectors.

Definition 3.3 (Selection (Yao et al., 2023)). A selection \oslash operates between two vectors $a \in \{0,1\}^d$, $b \in \mathbb{R}^d$ where $a \oslash b := [b_j : a_j = 1, j \in [d]]$.

Definition 3.4 (Invariant block selectors). The invariant block selectors $\Phi := \{\phi^{(i,k)}\}_{i \in [n_{\mathfrak{I}}], k \in V_i}$ with $\phi^{(i,k)} \in \{0,1\}^N$ perform selection (Defn. 3.3) on the encoded information: for any invariance property $\iota_i \in \mathfrak{I}$, any observable $\mathbf{x}^k, k \in V_i$ we have the selected representation:

$$\phi^{(i,k)} \oslash \hat{\mathbf{z}}^k = \phi^{(i,k)} \oslash g_k(\mathbf{x}^k) = \left[[g_k(\mathbf{x}^k)]_j : \phi_j^{(i,k)} = 1, j \in [N] \right],$$
(3.1)

with $\left\|\phi^{(i,k)}\right\|_0 = \|\phi^{(i,k')}\|_0 = |A_i|$ for all $\iota_i \in \mathfrak{I}, k, k' \in V_i$.

Intuition: Selectors select the relevant encoding dimensions for each invariance property $\iota_i \in \mathfrak{I}$. Each selector $\phi^{(i,k)}$ gives rise to a index subset $\hat{A}_i^k := \{j : \phi_j^{(i,k)} = 1\} \subseteq [N]$ that is specific to the invariance property ι_i and the observable \mathbf{x}^k . The assumption of known invariance size $|A_i|$ can be lifted in certain scenarios by, e.g., enforcing sharing between the learned latent variables, as shown by Fumero et al. (2024); Yao et al. (2023), or leveraging sparsity constraints (Lachapelle et al., 2022; 2024; Zheng et al., 2022; Xu et al., 2024).

Constraint 3.1 (Invariance constraint). For any invariance property $\iota_i \in \mathfrak{I}, i \in [n_{\mathfrak{I}}]$, the selected representations $\phi^{(i,k)} \oslash g_k(\mathbf{x}^k), k \in V_i$ must be ι_i -invariant across the observables from the subset $V_i \subseteq [K]$:

$$\iota_i(\phi^{(i,k)} \oslash g_k(\mathbf{x}^k)) = \iota_i(\phi^{(i,k')} \oslash g_{k'}(\mathbf{x}^{k'})) \quad \forall i \in [n_{\mathfrak{I}}] \; \forall k, k' \in V_i$$
(3.2)

Constraint 3.2 (Sufficiency constraint). For any $\iota_i \in \mathfrak{I}, i \in [n_{\mathfrak{I}}]$, the selected representation $\phi^{(i,k)} \oslash g_k(\mathbf{x}^k), k \in V_i$ must preserve all information of the invariant partition \mathbf{z}_{A_i} that we aim to identify, i.e., $I(\mathbf{z}_{A_i}, \phi^{(i,k)} \oslash g_k(\mathbf{x}^k)) = H(\mathbf{z}_{A_i}) \ \forall i \in [n_{\mathfrak{I}}], k \in V_i$, where $I(\cdot, \cdot)$ denotes the mutual information and $H(\cdot)$ denotes the differential entropy of the ground truth latent distribution $p_{\mathbf{z}_{A_i}}$.

300 **Remark:** The regularizer enforcing this sufficiency constraint can be tailored to suit the specific 301 task of interest. For example, for self-supervised training, it can be implemented as the mutual 302 information between the input data and the encodings, i.e., $I(\mathbf{x}, q(\mathbf{x})) = H(\mathbf{x})$, to preserve the 303 entropy from the observations; for classification, it becomes the mutual information between the 304 task labels and the learned representation $I(\mathbf{y}, g(\mathbf{x}))$. Sometimes, sufficiency does not have to be 305 enforced on the whole representation. For example, in the multiview line of work (von Kügelgen 306 et al., 2021; Daunhawer et al., 2023), when considering a single invariant block A, enforcing 307 sufficiency on the shared partition (implemented as entropy on the learned encoding $H(g(\mathbf{x})_{1:|A|})$) is enough to block-identify these shared latent variables z_A . 308

Theorem 3.1 (Identifiability of multiple invariant blocks). Consider a set of observables $S_{\mathbf{x}} = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^K\} \in \mathcal{X}$ generated from § 2 satisfying Asm. 2.1. Let G, Φ be the set of smooth encoders (Defn. 3.2) and selectors (Defn. 3.4) that satisfy Constraints 3.1 and 3.2, then the invariant component $\mathbf{z}_{A_i}^k$ is block-identified (Defn. 3.1) by $\phi^{(i,k)} \oslash g_k$ for all $\iota_i \in \mathfrak{I}, k \in [K]$.

314 **Discussion:** In general, Thm. 3.1 enforces all invariance properties $\iota_i \in \mathfrak{I}$ jointly and thus learns a representation that block-identifies all invariant blocks simultaneously. It allows mixing 315 multiple invariance principles, thus better adapting to complex real-world scenarios in which 316 various invariance relations typically occur. In practice, this constrained optimization problem 317 can be solved in many different flavors, e.g., Lippe et al. (2022b;a) employ a two-stage learning 318 process first to solve the sufficiency constraint and then the invariance constraint; Lachapelle 319 et al. (2023); Fumero et al. (2024) instead formulate it as a bi-level constrained optimization 320 problem. Some works (von Kügelgen et al., 2021; Daunhawer et al., 2023; Yao et al., 2023; 321 von Kügelgen et al., 2024; Zhang et al., 2024a; Ahuja et al., 2024) propose a loss that directly 322 solves the constrained optimization problem, while the others (Squires et al., 2023; Varici et al., 323 2024a;b) develop step-by-step algorithms as solutions.

What about the variant latents? Intuitively, the variant latents are not identifiable, as the invariance constraint (Constraint 3.1) is applied only to the selected invariant encodings, leaving the variant part without any weak supervision (Locatello et al., 2019). This result is formalized as follows:

Proposition 3.2 (General non-identifiability of variant latent variables). Consider the setup in Thm. 3.1, let $A := \bigcup_{i \in [n_{\mathfrak{I}}]} A_i$ denote the union of block-identified latent indices and $A^{c} := [N] \setminus A$ the complementary set where no ι -invariance $\iota \in \mathfrak{I}$ applies, then the variant latents $\mathbf{z}_{A^{c}}$ cannot be identified.

Although variant latent variables are generally non-identifiable, they can be identified under certain
 conditions. The following demonstrates that variant latent variables can be identified under invertible
 encoders when the variant and invariant partitions are mutually independent.

Proposition 3.3 (Identifiability of variant latent under independence). Consider an optimal encoder $g \in G^*$ and optimal selector $\phi \in \Phi^*$ from Thm. 3.1 that jointly identify an invariant block \mathbf{z}_A (we omit subscriptions k, i for simplicity), then $\mathbf{z}_{A^c}(A^c := [N] \setminus A)$ can be identified by the complementary encoding partition $(1 - \phi) \oslash g$ only if

- (i) g is invertible in the sense that $I(\mathbf{x}, g(\mathbf{x})) = H(\mathbf{x})$;
- (*ii*) \mathbf{z}_{A^c} *is independent on* \mathbf{z}_A .

Discussion: The generalization of new interventions has been a long-standing goal in causal representation learning. The generalization can be categorized into two layers: (1) generalize to unseen interventional values and (2) generalize to non-intervened nodes. The former includes the out-of-distributional value of the intervened node in the training set or a combination of multiple singly intervened nodes during training, which has been successfully demonstrated in various existing works (Zhang et al., 2024a; von Kügelgen et al., 2024). However, we argue that the second layer of generalization, namely generalizing to unseen nodes, is fundamentally impossible, as shown by Proposition 3.2; only under certain conditions such as independence and sufficient latent representation for reconstruction, non-intervened nodes in the training phase can be identified during inference (Proposition 3.3). This result aligns with the identifiability algebra given by (Yao et al., 2023) and is evidenced by numerous previous works, including disentanglement (Locatello et al., 2020; Fumero et al., 2022b; Lachapelle et al., 2022; 2024).

354 355 356

339

340

341 342

343

344

345

346

347

348

349

350

351

352

353

4 RELATED WORKS AS SPECIAL CASES OF OUR THEORY

This section provides an overview of the literature on causal representation learning (including multiview, multi-environment, temporal, and multi-task settings) and domain generalization, explaining the underlying invariance principles and data symmetries inherent in these works, which naturally fit into our framework as special cases. Tab. 1 provides a list of concrete examples and the explicit forms of their underlying invariance. Further mathematical details are deferred to App. D.

362 Multiview CRL. Multiview CRL (also termed "counterfactual" CRL) considers a setting where 363 each view (observable \mathbf{x}^k) is generated from a subset of latent causal variables (Locatello et al., 364 2020; Ahuja et al., 2022b; von Kügelgen et al., 2021; Daunhawer et al., 2023; Yao et al., 2023). 365 Given any set of jointly observed views, the view-specific generating latents could overlap, giving 366 rise to *sample level invariance* on all realizations of these shared latents. The common theoretical 367 contribution in this line of work in terms of identifiability is that the invariant partition of latents 368 (shared ones) can be block-identified by enforcing aligned and sufficient representation, which is a special case of Thm. 3.1 with specified sample invariance. 369

370 Multi-environment CRL. Multi-environment/interventional CRL (Ahuja et al., 2023; Squires 371 et al., 2023; Zhang et al., 2024a; Buchholz et al., 2024; Varici et al., 2023; 2024a; von Kügelgen 372 et al., 2021; Wendong et al., 2024) collects data from multiple environments that follow different 373 data distributions, often originated from interventions ($\mathbf{x}^k \sim P^k$). Current multi-environment CRL 374 literature has provided fruitful identifiability results based on various types of interventions: either 375 atomic or paired interventions per node or different parametric assumptions on the mixing function or the latent causal model. Interventions give rise to many types of invariance: When performing 376 an atomic intervention on an arbitrary node, the *marginal* of its non-descendants remain invariant; 377 the *score* of all other nodes than its parents and itself also remain invariant. By utilizing these two

378 types of invariance, we can not only explain various prior identification theories as special cases 379 of Thm. 3.1, but also directly develop new element-wise identification results on the latent variables, 380 given *imperfect* atomic interventions per node (Cor. D.1). Some other works (von Kügelgen et al., 381 2021; Varici et al., 2024a) consider paired interventions per node, with an *invariant interventional* 382 target between these paired interventional environments. This invariance imposes a certain score structure in the latent space, which can be used as an equivalent constraint as the invariant constraint (Constraint 3.1). More details in this regard are provided in App. D.2. More recently, 384 Ahuja et al. (2024) explains previous interventional identifiability results from a general weak distri-385 butional invariance perspective. Ahuja et al. (2024) proves block-affine identification (Defn. C.1) by 386 additionally assuming the mixing function to be finite degree polynomial, which can be explained 387 by Proposition C.2 together with our block-identifiability results under the general nonparametric 388 setting. They consider one single invariance set, which is a special case of Thm. 3.1 with one joint 389 ι -property. Another line of interventional CRL work (Zhang et al., 2024a) employs an orthogonal 390 proof technique, originating from nonlinear ICA with auxiliary variables (Hyvarinen et al., 2019). 391 We remark that our framework does not directly include this line of identifiability theory.

392

404

393 Temporal CRL. Extending causal representation learning into time-series setting, temporal CRL often assumes an "intervenable" trajectory in the latent space (Lippe et al., 2022a;b; 2023; Lachapelle 394 et al., 2022; Yao et al., 2022b;a; Li et al., 2024b). At each time step, an intervention/action modifies 395 the dynamics of a subset of latent variables, with the remaining invariant partition following the 396 default dynamics conditioning on the previous time step. Existing works have shown that the 397 intervened part can be disentangled from the invariant part when there is no causal link between 398 the latent causal variables at the same time step (Lachapelle et al., 2022; Lippe et al., 2022a;b). 399 Comparing the "counterfactual" latent with the actual partially intervened latents on the same time 400 step, one observes the *transitional distribution* (current latents conditioning on previous latents) 401 remain invariant for the non-intervened partition. This formulates an explicit ι -property (Defn. 2.1) 402 for each time step with potentially different invariant partitions, explaining many existing temporal 403 CRL identifiability theories by incorporating Thm. 3.1.

Multi-task CRL In supervised CRL, latent variables (Lachapelle et al., 2023; Fumero et al., 405 2024) are shown to be identifiable under multi-task setting, meaning there are multiple task labels 406 available for each observable ($\mathbf{x}^k := (\mathbf{x}, \mathbf{y}^k)$). The key criterion for achieving identifiability is 407 overlapping task support, i.e., a set of tasks depends on a shared set of latents. Defining this shared 408 set of latent as the invariant partition z_A in our setup, we obtain a valid *i*-property (Defn. 2.1) 409 defined by the optimal classifier of individual tasks (Details provided in App. D.4). Incorporating 410 this invariance principle into Thm. 3.1 explains the identification results of (Lachapelle et al., 2023; 411 Fumero et al., 2024), showing the overlapping task support can be identified. 412

413 Domain Generalization. The field of domain generalization focuses on the out-of-distribution 414 performance of the learned representation instead of the theoretical identifiability guarantee (Rojas-Carulla et al., 2018; Arjovsky et al., 2020; Ahuja et al., 2022a; Krueger et al., 2021; Sagawa et al., 415 2019). The goal is to learn representations that perform equally well across domains originating 416 from distributional shifts, such as covariates shift or concept shift. Domain generalization typically 417 assumes the same downstream prediction task, and this task depends on the same subset of latent 418 factors A across all domains. Given the same ground truth task-latent dependency, the *domain risk* 419 w.r.t. ground truth inverting process remains invariant across all domains. This invariance property 420 together with Thm. 3.1 could provide theoretical insights for domain generalization works such 421 as (Krueger et al., 2021; Sagawa et al., 2019) (formal mathematical derivation provided in (f)).

422 423 424

5 EXPERIMENTS

This section demonstrates the real-world applicability of causal representation learning under the invariance principle, evidenced by superior treatment effect estimation performance on the high-dimensional causal inference benchmark (Cadei et al., 2024) using a regularizer for the domain generalization literature that utilizes the invariance principle (Krueger et al., 2021) (§ 5.1). Additionally, we provide ablation studies on existing interventional causal representation learning methods (Wendong et al., 2024; Ahuja et al., 2023; von Kügelgen et al., 2024), showcasing that non-trivial distributional invariance is needed for latent variable identification. This distributional invariance could, but does not have to, arise from a valid intervention in the sense of causality (§ 5.2).



Figure 1: TERB and Balanced Accuracy with standard deviation over 20 different seeds varying the invariance weight λ_{INV} of V-REx (Krueger et al., 2021) on ISTAnt dataset (Cadei et al., 2024). Stars represent the selected best models based on a small but heterogeneous validation set.

449 5.1 CASE STUDY: ISTANT

445

446

447 448

460 461 462

This experiment focuses on ISTAnt (Cadei et al., 2024), a recent real-world ecological benchmark designed for treatment effect estimation. ISTAnt consists of video recordings of ants triplets with occasional grooming behavior. The goal is to extract a per-frame representation for supervised behavior
classification (grooming or not) to estimate the Average Treatment Effect of an intervention (exposure to a chemical substance). Further details about the problem setting are provided in App. F.1.

Experiment settings. Different videos in ISTAnt are considered different *experiments* as the experiment settings and treatments vary. We consider hard annotation sampling criteria (more non-annotated than annotated) for both experiments (videos) and positions, as described by Cadei et al. (2024). For the training, we adopt a domain generalization objective that utilizes the invariance principle (Krueger et al., 2021), which is restated as follows:

$$\mathcal{R}_{\text{V-REx}}(\mathbf{w} \circ g) = \underbrace{\lambda_{\text{INV}} \operatorname{Var}(\{\mathcal{R}_1(\mathbf{w} \circ g), \dots, \mathcal{R}_K(\mathbf{w} \circ g)\})}_{\text{invariance}} + \underbrace{\sum_{k \in [K]} \mathcal{R}_k(\mathbf{w} \circ g)}_{\text{sufficiency}}, \tag{5.1}$$

we provide a detailed derivation in (f) showing the invariance term above is indeed enforcing risk invariance. We vary the strength of the invariant component in eq. (5.1) by setting the regularization multiplier λ_{INV} from 0 (ERM) to 10 000. We repeat 20 independent runs for each λ_{INV} to estimate the statistical error. All other implementational details follow Cadei et al. (2024). We evaluate the performance with both *balanced accuracy* and *Treatment Effect Relative Bias* (TERB). TERB is defined by Cadei et al. (2024) as the ratio between the bias in the predictions across treatment groups and the true average treatment effect estimated with ground-truth annotations over the whole trial.

470 **Results.** Fig. 1 depicts the model performance regarding varying invariance regularization strength 471 $\lambda_{\rm INV}$. As expected, the balanced accuracy initially increases with the $\lambda_{\rm INV}$, as adequate invariance 472 enforces identifying task-related latents, thus benefiting the prediction problem. At a later point, 473 the performance decreases because the sufficiency component is not correctly balanced with the 474 invariance. Similarly, the TERB improves positively, weighting the invariance component until a certain threshold. On average, with $\lambda_{INV} = 100$ the TERB decreases to 20% (from 100% using 475 ERM) with experiment subsampling. In agreement with (Cadei et al., 2024), a naive estimate of the 476 TEB on a small validation set is a reasonable (albeit not perfect) model selection criterion. Although 477 it performs slightly worse than model selection based on ERM loss in the position sampling case, 478 it shows more reliability overall. This experiment underscores the advantages of flexibly enforcing 479 known invariances in the data, corroborating our identifiability theory $(\S 3)$. 480

481 5.2 Synthetic Ablation with "Ninterventions"

This subsection presents identifiability results under controversial (non-causal) conditions using simulated data. We consider a simple graph of three causal variables as $\mathbf{z}_1 \rightarrow \mathbf{z}_2 \rightarrow \mathbf{z}_3$. The corresponding joint density has the form of

$$p_{\mathbf{z}}(z_1, z_2, z_3) = p(z_3 \mid z_2)p(z_2 \mid z_1)p(z_1)$$

This experiment aims to demonstrate that existing methods for interventional CRL rely primarily
 on distributional invariance, regardless of whether this invariance arises from a well-defined intervention or some other arbitrary transformation. To illustrate this, we introduce the concept of a
 "inintervention," which has a similar distributional effect to a regular intervention, maintaining certain conditionals invariant while altering others, but without a causal interpretation.

Definition 5.1 (Nintervention). We define a "*nintervention*" on a causal conditional as the process of changing its distribution but cutting all incoming and outgoing edges. Child nodes condition on the old, pre-intervention, random variable. Formally, we consider the latent SCM as defined in Defn. B.1, an *nintervention* on a node $j \in [N]$ gives rise to the following conditional factorization

$$\tilde{p}_{\mathbf{z}}(z) = \tilde{p}(z_j) \prod_{i \in [N] \setminus \{j\}} p(z_i \mid z_{\text{pa}(i)}^{\text{old}})$$

Note that the marginal distribution of all non-nintervened nodes $P_{\mathbf{z}_{[N]\setminus j}}$ remain invariant after nintervention. In previous example, we perform a nintervention by replacing the conditional density $p(z_2 | z_1)$ using a sufficiently different marginal distribution $p(\tilde{z}_2)$ that satisfies Defn. 2.1 (ii), which gives rise to the following new factorization $\tilde{p}_{\mathbf{z}}(z_1, z_2, z_3) = p(z_3 | z_2^{\text{old}})\tilde{p}(z_2)p(z_1)$. Note that \mathbf{z}_3 conditions on the random variable \mathbf{z}_2 before nintervention, whose realization is denoted as z_2^{old} . Differing from a causal *intervention*, we cut both the incoming and outgoing links of \mathbf{z}_2 and keep the marginal distribution of \mathbf{z}_3 the same. Clearly, this is a non-sensical intervention from the causal perspective because we eliminate the causal effect from \mathbf{z}_2 to its descendants.

506 Experiment settings. As a proof of concept, we choose a linear Gaussian additive noise model and 507 a nonlinear mixing function implemented as a 3-layer invertible MLP. We average the results over 508 three independently sampled *ninterventional* densities $\tilde{p}(z_2)$ while guaranteeing all *ninterventional* 509 distributions satisfy Defn. 2.1 (ii). As the marginal distribution of both z_1, z_3 remains the same after 510 a *nintervention*, we expect z_1, z_3 to be block-identified (Defn. 3.1) according to Thm. 3.1. In prac-511 tice, we enforce the marginal invariance constraint (Constraint 3.1) by minimizing the MMD loss, as 512 implemented by the interventional CRL works (Zhang et al., 2024a; Ahuja et al., 2024) and train an 513 auto-encoder for a sufficient representation (Constraint 3.2). Further details are included in App. F.

Results. To validate block-identifiability, we perform Kernel-Ridge Regression between the estimated block $[\hat{z}1, \hat{z}3]$ and the ground truth latents z_1, z_2, z_3 . Both z_1 and z_3 are block-identified with high R^2 scores of 0.863 ± 0.031 and 0.872 ± 0.035 . In contrast, z_2 is not identified, with a low R^2 of 0.065 ± 0.017 , indicating identification is driven by the underlying distributional invariance.

519 520 6 CONCLUSIONS

491

492

493

494

495 496

497

521 In this paper, we take a closer look at the wide range of causal representation learning methods. Interestingly, we find many CRL approaches share methodological similarities in aligning the 522 representation to known data symmetries. We identified two components involved in identifiability 523 results: preserving information of the data and a set of known invariances (§ 3). Our results help 524 clarify the role of causal assumptions in causal variable identification, shifting the focus from a 525 characterization of specific assumptions for identifiability, which are not necessarily satisfied in 526 real-world scenarios, to a general recipe that allows practitioners to specify known invariances in 527 their problem and learn representations that align with them. Following the general recipe, we 528 successfully exemplified the real-world applicability of CRL on ecological data, as shown in § 5.1. 529 Nevertheless, our paper leaves out certain settings concerning identifiability that may be interesting 530 for future work, such as discrete variables and finite sample guarantees. 531

532 ETHICS STATEMENT

This work contributes to causal representation learning by unifying many existing theoretical results,
thus vastly broadening its real-world applicability. As the paper is predominantly theoretical, we
believe it poses no immediate ethical risks.

536 537 Reproducibility Statement

All proofs in this paper are deferred to App. E. The ISTAnt dataset in § 5.1 is published by (Cadei et al., 2024). Results provided in § 5 can be reproduced following the details given in App. F. The curated code is included in the supplementary material and will be published upon acceptance.

540 REFERENCES 541

552

553

554

555

558

565

566

567

568

569

570

576

581

- Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Jean-Christophe Gagnon-Audet, Yoshua Bengio, 542 Ioannis Mitliagkas, and Irina Rish. Invariance principle meets information bottleneck for out-of-543 distribution generalization, 2022a. 1, 2, 8, 27, 44 544
- Kartik Ahuja, Jason S Hartford, and Yoshua Bengio. Weakly supervised representation learning 546 with sparse perturbations. Advances in Neural Information Processing Systems, 35:15516–15528, 547 2022b. 5, 7, 22, 23, 42 548
- 549 Kartik Ahuja, Divyat Mahajan, Yixin Wang, and Yoshua Bengio. Interventional causal representa-550 tion learning. In International Conference on Machine Learning, pp. 372–407. PMLR, 2023. 1, 7, 8, 20, 24, 29, 33 551
 - Kartik Ahuja, Amin Mansouri, and Yixin Wang. Multi-domain causal representation learning via weak distributional invariances. In International Conference on Artificial Intelligence and Statistics, pp. 865-873. PMLR, 2024. 6, 8, 10, 23, 25, 39
- 556 Joshua D Angrist and Jörn-Steffen Pischke. Mostly harmless econometrics: An empiricist's com*panion*. Princeton university press, 2009. 1
- J Antonakis and R Lalive. Counterfactuals and causal inference: Methods and principles for social 559 research. Structural Equation Modeling, 18(1):152–159, 2011. 1 560
- 561 Md Rifat Arefin, Yan Zhang, Aristide Baratin, Francesco Locatello, Irina Rish, Dianbo Liu, and 562 Kenji Kawaguchi. Unsupervised concept discovery mitigates spurious correlations. In *Forty-first* 563 International Conference on Machine Learning, 2024. 37
 - Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization, 2020. 2, 8, 27, 44
 - Jinze Bai, Rui Men, Hao Yang, Xuancheng Ren, Kai Dang, Yichang Zhang, Xiaohuan Zhou, Peng Wang, Sinan Tan, An Yang, et al. Ofasys: A multi-modal multi-task learning system for building generalist models. arXiv preprint arXiv:2212.04408, 2022. 26
- 571 Elias Bareinboim and Judea Pearl. Causal inference and the data-fusion problem. Proceedings of 572 the National Academy of Sciences, 113(27):7345-7352, 2016. 1 573
- 574 Sander Beckers and Joseph Y Halpern. Abstracting causal models. In Proceedings of the aaai conference on artificial intelligence, volume 33, pp. 2678–2685, 2019. 22 575
- Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In Proceedings of 577 the European conference on computer vision (ECCV), pp. 456–473, 2018. 27 578
- 579 Gilles Blanchard, Gyemin Lee, and Clayton Scott. Generalizing from several related classification 580 tasks to a new unlabeled sample. Advances in neural information processing systems, 24, 2011. 27 582
- Johann Brehmer, Pim De Haan, Phillip Lippe, and Taco S Cohen. Weakly supervised causal repre-583 sentation learning. Advances in Neural Information Processing Systems, 35:38319–38331, 2022. 584 1, 4, 21, 22, 42 585
- 586 Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geomet-587 ric deep learning: going beyond euclidean data. IEEE Signal Processing Magazine, 34(4):18–42, 588 2017. 2. 37 589
- Michael M Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning: 591 Grids, groups, graphs, geodesics, and gauges. arXiv preprint arXiv:2104.13478, 2021. 2, 37
- Glen D Brown, Satoshi Yamada, and Terrence J Sejnowski. Independent component analysis at the 593 neural cocktail party. Trends in neurosciences, 24(1):54-63, 2001. 1

606

607

608

613

619

625

626

627

628

633

394	Simon Buchholz, Goutham Rajendran, Elan Rosenfeld, Bryon Aragam, Bernhard Schölkopf, and
595	Pradeep Ravikumar. Learning linear causal representations from interventions under general non-
596	linear mixing. Advances in Neural Information Processing Systems, 36, 2024. 1, 7, 19, 23, 24,
597	28, 39

- Riccardo Cadei, Lukas Lindorfer, Sylvia Cremer, Cordelia Schmid, and Francesco Locatello. Smoke
 and mirrors in causal downstream tasks. *Advances in Neural Information Processing Systems*, 37, 2024. 1, 2, 8, 9, 10, 34, 35
- Rich Caruana. Multitask learning. *Machine learning*, 28:41–75, 1997. 26
- Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International confer- ence on machine learning*, pp. 2990–2999. PMLR, 2016. 2, 38
 - Taco S Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Spherical cnns. *arXiv preprint arXiv:1801.10130*, 2018. 38
- Elliot Creager, Joern-Henrik Jacobsen, and Richard Zemel. Environment inference for invariant
 learning. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Con- ference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp.
 2189–2200. PMLR, 18–24 Jul 2021. 37
- Imant Daunhawer, Alice Bizeul, Emanuele Palumbo, Alexander Marx, and Julia E Vogt. Identifia bility results for multimodal contrastive learning. In *The Eleventh International Conference on Learning Representations*, 2023. 4, 5, 6, 7, 22, 23, 42
- Thomas Dean and Keiji Kanazawa. A model for reasoning about persistance and causation. In
 Computational Intelligence, pp. 5, 1989. 25
- David L Donoho. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–1306, 2006. 1
- Chen Fang, Ye Xu, and Daniel N Rockmore. Unbiased metric learning: On the utilization of multiple
 datasets and web images for softening bias. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1657–1664, 2013. 27
 - Marco Fumero, Luca Cosmo, Simone Melzi, and Emanuele Rodolà. Learning disentangled representations via product manifold projection. In *International conference on machine learning*, pp. 3530–3540. PMLR, 2021. 38
- Marco Fumero, Florian Wenzel, Luca Zancato, Alessandro Achille, Emanuele Rodolà, Stefano
 Soatto, Bernhard Schölkopf, and Francesco Locatello. Leveraging sparse and shared feature
 activations for disentangled representation learning. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 4, 5, 6, 7, 8, 26, 27, 28, 43
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of machine learning research*, 17(59):1–35, 2016. 27
- Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. arXiv preprint arXiv:2007.01434, 2020. 35
- Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2):20170016, 2018. 2
- Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende, and Alexander Lerchner. Towards a definition of disentangled representations. *arXiv preprint arXiv:1812.02230*, 2018. 38
- Patrik Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear
 causal discovery with additive noise models. *Advances in neural information processing systems*, 21, 2008. 21

- 648 Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. The platonic representation 649 hypothesis. arXiv preprint arXiv:2405.07987, 2024. 37 650 Aapo Hyvarinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ica using auxiliary variables and 651 generalized contrastive learning. In The 22nd International Conference on Artificial Intelligence 652 and Statistics, pp. 859-868. PMLR, 2019. 8, 29 653 654 Dominik Janzing and Sergio Hernan Garrido Mejia. A phenomenological account for causality in 655 terms of elementary actions. Journal of Causal Inference, 12(1):20220076, 2024. 37 656 657 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, 658 Kathryn Tunyasuyunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate 659 protein structure prediction with alphafold. nature, 596(7873):583–589, 2021. 38 660 Ilyes Khemakhem, Ricardo Monti, Diederik Kingma, and Aapo Hyvarinen. Ice-beem: Identifiable 661 conditional energy-based deep models based on nonlinear ica. Advances in Neural Information 662 Processing Systems, 33:12768–12778, 2020. 26, 29 663 664 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint 665 arXiv:1312.6114, 2013. 26 666 667 Bohdan Kivva, Goutham Rajendran, Pradeep Ravikumar, and Bryon Aragam. Identifiability of deep generative models without auxiliary information. In S. Koyejo, S. Mohamed, A. Agarwal, 668 D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems, 669 volume 35, pp. 15687-15701. Curran Associates, Inc., 2022. 29 670 671 David A Klindt, Lukas Schott, Yash Sharma, Ivan Ustyuzhaninov, Wieland Brendel, Matthias 672 Bethge, and Dylan Paiton. Towards nonlinear disentanglement in natural data with temporal 673 sparse coding. In International Conference on Learning Representations, 2021. 26 674 Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural 675 network representations revisited. In International conference on machine learning, pp. 3519-676 3529. PMLR, 2019. 37 677 678 David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai 679 Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapola-680 tion (rex). In International conference on machine learning, pp. 5815–5826. PMLR, 2021. 1, 2, 681 8, 9, 27, 29, 35, 37, 44 682 Sébastien Lachapelle, Rodriguez Lopez, Pau, Yash Sharma, Katie E. Everett, Rémi Le Priol, Alexan-683 dre Lacoste, and Simon Lacoste-Julien. Disentanglement via mechanism sparsity regularization: 684 A new principle for nonlinear ICA. In First Conference on Causal Learning and Reasoning, 2022. 685 1, 4, 5, 6, 7, 8, 19, 20, 25, 26, 29 686 687 Sébastien Lachapelle, Tristan Deleu, Divyat Mahajan, Ioannis Mitliagkas, Yoshua Bengio, Simon 688 Lacoste-Julien, and Quentin Bertrand. Synergies between disentanglement and sparsity: Gen-689 eralization and identifiability in multi-task learning. In International Conference on Machine 690 Learning, pp. 18171–18206. PMLR, 2023. 1, 4, 5, 6, 8, 26, 28, 29, 43 691 Sébastien Lachapelle, Pau Rodríguez López, Yash Sharma, Katie Everett, Rémi Le Priol, Alexan-692 dre Lacoste, and Simon Lacoste-Julien. Nonparametric partial disentanglement via mecha-693 nism sparsity: Sparse actions, interventions and sparse temporal dependencies. arXiv preprint 694 arXiv:2401.04890, 2024. 6, 7, 25, 26, 29 695 696 Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain 697 generalization. In Proceedings of the IEEE international conference on computer vision, pp. 698 5542-5550, 2017. 27 699 Yixuan Li, Jason Yosinski, Jeff Clune, Hod Lipson, and John Hopcroft. Convergent learning: Do 700 different neural networks learn the same representations? arXiv preprint arXiv:1511.07543, 2015. 701

702 703 704 705	Zijian Li, Ruichu Cai, Zhenhui Yang, Haiqin Huang, Guangyi Chen, Yifan Shen, Zhengming Chen, Xiangchen Song, Zhifeng Hao, and Kun Zhang. When and how: Learning identifiable latent states for nonstationary time series forecasting. <i>arXiv preprint arXiv:2402.12767</i> , 2024a. 25
706 707 708	Zijian Li, Yifan Shen, Kaitao Zheng, Ruichu Cai, Xiangchen Song, Mingming Gong, Zhifeng Hao, Zhengmao Zhu, Guangyi Chen, and Kun Zhang. On the identification of temporally causal representation with instantaneous dependence. <i>arXiv preprint arXiv:2405.15325</i> , 2024b. 8, 25, 26
709 710 711 712	Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M Asano, Taco Cohen, and Efstratios Gavves. Causal representation learning for instantaneous and temporal effects in interactive systems. In <i>The Eleventh International Conference on Learning Representations</i> , 2022a. 1, 4, 5, 6, 8, 25, 26, 43
713 714 715 716	Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M Asano, Taco Cohen, and Stratis Gavves. Citris: Causal identifiability from temporal intervened sequences. In <i>International Conference on Machine Learning</i> , pp. 13557–13603. PMLR, 2022b. 1, 4, 5, 6, 7, 8, 25, 26, 30, 42
717 718 719	Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M Asano, Taco Cohen, and Efstratios Gavves. Biscuit: Causal representation learning from binary interactions. In Uncertainty in Artificial Intelligence, pp. 1263–1273. PMLR, 2023. 1, 8, 25, 26, 43
720 721 722 723	Yuejiang Liu, Alexandre Alahi, Chris Russell, Max Horn, Dominik Zietlow, Bernhard Schölkopf, and Francesco Locatello. Causal triplet: An open challenge for intervention-centric causal repre- sentation learning. In <i>Conference on Causal Learning and Reasoning</i> , pp. 553–573. PMLR, 2023. 22
724 725 726 727 728	Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In <i>international conference on machine learning</i> , pp. 4114–4124. PMLR, 2019. 7
729 730 731 732 733	Francesco Locatello, Ben Poole, Gunnar Raetsch, Bernhard Schölkopf, Olivier Bachem, and Michael Tschannen. Weakly-supervised disentanglement without compromises. In Hal Daumé III and Aarti Singh (eds.), Proceedings of the 37th International Conference on Machine Learn- ing, volume 119 of Proceedings of Machine Learning Research, pp. 6348–6359. PMLR, 13–18 Jul 2020. 5, 7, 22, 23, 42
734 735 736 737	Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello, and Emanuele Rodolà. Relative representations enable zero-shot latent space communication. <i>Inter-</i> <i>national Conference on Learning Representations</i> , 2022. 37
738 739	Kevin Patrick Murphy. <i>Dynamic bayesian networks: representation, inference and learning</i> . University of California, Berkeley, 2002. 25
740 741 742 743	Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. <i>arXiv preprint arXiv:2304.07193</i> , 2023. 35
744	Judea Pearl. Causality. Cambridge university press, 2009. 21
745 746 747 748	Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pp. 1406–1415, 2019. 27
749 750	Jonas Peters, Joris M Mooij, Dominik Janzing, and Bernhard Schölkopf. Causal discovery with continuous additive noise models. <i>Journal of Machine Learning Research</i> , 2014. 2, 37
751 752 753	Jonas Peters, Dominik Janzing, and Bernhard Schlkopf. <i>Elements of Causal Inference: Foundations and Learning Algorithms</i> . The MIT Press, 2017. ISBN 0262037319. 37
754 755	Mohammad Pezeshki, Diane Bouchacourt, Mark Ibrahim, Nicolas Ballas, Pascal Vincent, and David Lopez-Paz. Discovering environments with xrm. In <i>Forty-first International Conference on Machine Learning</i> , 2024. 37

756 757 758	Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In <i>Interna-</i> <i>tional conference on machine learning</i> , pp. 1530–1538. PMLR, 2015. 26
759 760 761	Tapani Ristaniemi. On the performance of blind source separation in cdma downlink. In <i>Proceedings of the International Workshop on Independent Component Analysis and Signal Separation (ICA'99)</i> , pp. 437–441, 1999. 1
762 763 764	James M Robins, Miguel Angel Hernan, and Babette Brumback. Marginal structural models and causal inference in epidemiology. <i>Epidemiology</i> , 11(5):550–560, 2000. 35
765 766	Mateo Rojas-Carulla, Bernhard Schölkopf, Richard Turner, and Jonas Peters. Invariant models for causal transfer learning. <i>Journal of Machine Learning Research</i> , 19(36):1–34, 2018. 1, 8
767 768 769 770	P Rubenstein, S Weichwald, S Bongers, J Mooij, D Janzing, M Grosse-Wentrup, and B Schölkopf. Causal consistency of structural equation models. In <i>33rd Conference on Uncertainty in Artificial</i> <i>Intelligence (UAI 2017)</i> , pp. 808–817. Curran Associates, Inc., 2017. 22
771 772	Donald B Rubin. Causal inference using potential outcomes: Design, modeling, decisions. <i>Journal of the American Statistical Association</i> , 100(469):322–331, 2005. 22
773 774 775	Jakob Runge. Modern causal inference approaches to investigate biodiversity-ecosystem functioning relationships. <i>nature communications</i> , 14(1):1917, 2023. 35
776 777 778	Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. <i>arXiv preprint arXiv:1911.08731</i> , 2019. 1, 2, 8, 27, 28, 29, 43
779 780 781 782	Jonathan M Samet, Francesca Dominici, Frank C Curriero, Ivan Coursac, and Scott L Zeger. Fine particulate air pollution and mortality in 20 us cities, 1987–1994. <i>New England journal of medicine</i> , 343(24):1742–1749, 2000. 35
783 784 785	Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. <i>Proceedings of the IEEE</i> , 109(5):612–634, 2021. 1, 29
786 787 788	Maximilian Seitzer, Bernhard Schölkopf, and Georg Martius. Causal influence detection for improv- ing efficiency in reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 34: 22905–22918, 2021. 1
789 790 791 792	Chandler Squires, Anna Seigal, Salil S. Bhate, and Caroline Uhler. Linear causal disentanglement via interventions. In <i>International Conference on Machine Learning</i> , volume 202, pp. 32540–32560. PMLR, 2023. 5, 6, 7, 21, 23, 24, 28, 39
793 794 795	Núria Armengol Urpí, Marco Bagatella, Marin Vlastelica, and Georg Martius. Causal action influ- ence aware counterfactual data augmentation. <i>Forty-first International Conference on Machine</i> <i>Learning</i> , 2024. 1
796 797 798 799	Egbert H Van Nes, Marten Scheffer, Victor Brovkin, Timothy M Lenton, Hao Ye, Ethan Deyle, and George Sugihara. Causal feedbacks in climate change. <i>Nature Climate Change</i> , 5(5):445–448, 2015. 35
800 801 802	 Burak Varici, Emre Acarturk, Karthikeyan Shanmugam, Abhishek Kumar, and Ali Tajer. Score-based causal representation learning with interventions. <i>arXiv preprint arXiv:2301.08230</i>, 2023. 4, 7, 23, 24, 25, 28, 30, 34, 40
803 804 805 806 807	Burak Varici, Emre Acartürk, Karthikeyan Shanmugam, and Ali Tajer. General identifiability and achievability for causal representation learning. In <i>International Conference on Artificial Intelligence and Statistics</i> , pp. 2314–2322. PMLR, 2024a. 1, 4, 6, 7, 8, 20, 21, 23, 24, 25, 28, 30, 40
808 809	Burak Varici, Emre Acartürk, Karthikeyan Shanmugam, and Ali Tajer. Linear causal representation learning from unknown multi-node interventions. <i>arXiv preprint arXiv:2406.05937</i> , 2024b. 4, 5, 6, 20, 23, 24, 28, 30, 40, 41

810 Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep 811 hashing network for unsupervised domain adaptation. In Proceedings of the IEEE conference on 812 computer vision and pattern recognition, pp. 5018-5027, 2017. 27 813 Ricardo Vigário, Veikko Jousmäki, Matti Hämäläinen, Riitta Hari, and Erkki Oja. Independent com-814 ponent analysis for identification of artifacts in magnetoencephalographic recordings. Advances 815 in neural information processing systems, 10, 1997. 1 816 817 Julius von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel 818 Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably 819 isolates content from style. Advances in neural information processing systems, 34:16451–16467, 820 2021. 4, 5, 6, 7, 8, 20, 22, 23, 24, 34, 41 821 Julius von Kügelgen, Michel Besserve, Liang Wendong, Luigi Gresele, Armin Kekić, Elias Barein-822 boim, David Blei, and Bernhard Schölkopf. Nonparametric identifiability of causal represen-823 tations from unknown interventions. Advances in Neural Information Processing Systems, 36, 824 2024. 1, 4, 6, 7, 8, 19, 20, 21, 24, 25, 28, 30, 41 825 826 Liang Wendong, Armin Kekić, Julius von Kügelgen, Simon Buchholz, Michel Besserve, Luigi Gre-827 sele, and Bernhard Schölkopf. Causal component analysis. Advances in Neural Information 828 Processing Systems, 36, 2024. 4, 7, 8, 28, 30, 41 829 Danru Xu, Dingling Yao, Sébastien Lachapelle, Perouz Taslakian, Julius von Kügelgen, Francesco 830 Locatello, and Sara Magliacane. A sparsity principle for partially observable causal representation 831 learning. Forty-first International Conference on Machine Learning, 2024. 6, 29 832 833 Dingling Yao, Danru Xu, Sebastien Lachapelle, Sara Magliacane, Perouz Taslakian, Georg Martius, 834 Julius von Kügelgen, and Francesco Locatello. Multi-view causal representation learning with partial observability. In The Twelfth International Conference on Learning Representations, 2023. 835 4, 5, 6, 7, 20, 22, 23, 42 836 837 Dingling Yao, Caroline Muller, and Francesco Locatello. Marrying causal representation learning 838 with dynamical systems for science. Advances in Neural Information Processing Systems, 37, 839 2024. 36 840 841 Weiran Yao, Guangyi Chen, and Kun Zhang. Temporally disentangled representation learning. Advances in Neural Information Processing Systems, 35:26492–26503, 2022a. 8, 25, 26 842 843 Weiran Yao, Yuewen Sun, Alex Ho, Changyin Sun, and Kun Zhang. Learning temporally causal 844 latent processes from general temporal data. International Conference on Learning Representa-845 tions, 2022b. 8, 25, 26 846 847 Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical 848 risk minimization. arXiv preprint arXiv:1710.09412, 2017. 27 849 Jiaqi Zhang, Kristjan Greenewald, Chandler Squires, Akash Srivastava, Karthikeyan Shanmugam, 850 and Caroline Uhler. Identifiability guarantees for causal disentanglement from soft interventions. 851 Advances in Neural Information Processing Systems, 36, 2024a. 1, 5, 6, 7, 8, 10, 20, 21, 23, 24, 852 25, 28, 33, 41 853 854 K Zhang and A Hyvärinen. On the identifiability of the post-nonlinear causal model. In 25th 855 Conference on Uncertainty in Artificial Intelligence (UAI 2009), pp. 647–655. AUAI Press, 2009. 21.22 856 857 Kun Zhang and Aapo Hyvärinen. Distinguishing causes from effects using nonlinear acyclic causal 858 models. In Causality: Objectives and Assessment, pp. 157-164. PMLR, 2010. 21 859 Kun Zhang, Shaoan Xie, Ignavier Ng, and Yujia Zheng. Causal representation learning from multi-861 ple distributions: A general setting. Internatinal Conference on Machine Learning, 2024b. 29 862 Yu Zhang and Qiang Yang. An overview of multi-task learning. National Science Review, 5(1): 863 30-43, 2018. 26

864 865 866	Yuli Zhang, Huaiyu Wu, and Lei Cheng. Some new deformation formulas about variance and covariance. In 2012 proceedings of international conference on modelling, identification and control, pp. 987–992. IEEE, 2012. 29
868 869	Yujia Zheng, Ignavier Ng, and Kun Zhang. On the identifiability of nonlinear ica: Sparsity and beyond. <i>Advances in neural information processing systems</i> , 35:16411–16422, 2022. 6, 29
870 871 872	Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A survey. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 45(4):4396–4415, 2022. 27
873 874 875 876	Jinguo Zhu, Xizhou Zhu, Wenhai Wang, Xiaohua Wang, Hongsheng Li, Xiaogang Wang, and Jifeng Dai. Uni-perceiver-moe: Learning sparse generalist models with conditional moes. <i>Advances in Neural Information Processing Systems</i> , 35:2664–2678, 2022. 26
877 878 879	Roland S Zimmermann, Yash Sharma, Steffen Schneider, Matthias Bethge, and Wieland Brendel. Contrastive learning inverts the data generating process. In <i>International Conference on Machine Learning</i> , pp. 12979–12990. PMLR, 2021. 5
880 881	
882	
883	
884	
885	
886	
887	
888	
889	
890	
891	
892	
893	
894	
895	
090 207	
898	
899	
900	
901	
902	
903	
904	
905	
906	
907	
908	
909	
910	
911	
912	
913	
914	
916	
917	

Appendix

Table of Contents

Α	Notation and Terminology	18
В	Preliminaries	19
С	Identifiability Theory C.1 On the granularity of identification C.2 Identifying the causal graph	19 19 21
D	Related WorksD.1Multiview Causal Representation LearningD.2Multi-environment Causal Representation LearningD.3Temporal Causal Representation LearningD.4Multi-task Causal Representation LearningD.5Domain GeneralizationD.6Further Explanations for Tab. 4D.7Notable Cases Not Directly Covered by the Theory	22 22 23 25 26 27 28 29
Ε	ProofsE.1Assumption JustificationE.2Proof for Thm. 3.1E.3Proofs for Generalization of Variant LatentsE.4Proofs for Granularity of Latent Variable IdentificationE.5Proof for Cor. D.1	29 29 30 31 33 34
F	Implementation DetailsF.1Case Study: ISTAntF.2Synthetic Ablation with "Ninterventions"	34 34 36
G	Further Discussions and Connections to Other Fields G.1 Representational Alignment and Platonic Representation G.2 Environment Discovery G.3 Geometric Deep Learning	36 37 37 37

A NOTATION AND TERMINOLOGY

956		
957	f	Mixing function
958 959	g	Smooth encoder
960	${\cal G}$	Ground truth causal graph
961 962	x	Entangled observables
963	\mathbf{Z}	Ground truth latent variables
964	D	Dimensionality of observable x
965 966	N	Dimensionality of latents z
967	A	Subset of latent indices with invariance properties $(A \subseteq [N])$
968 969	ι	Projector which maps the latents to the space where the invariance property holds
970	\sim_{ι}	The latent equivalence relation
971	I	A set of invariance properties

- 972 \mathcal{X} Support of a set of observables $\mathcal{S}_{\mathbf{x}}$ 973
- 974 Z Support of a set of latent vectors S_z
- 975 G A set of smooth encoders
- 976 977 Φ A set of selectors
- 978 TC Transitive closure

980 B PRELIMINARIES

979

996 997 998

1024

In this subsection, we revisit the common definitions and assumptions in identifiability works from
 causal representation learning. We begin with the definition of a latent structural causal model:

Definition B.1 (Latent SCM (von Kügelgen et al., 2024)). Let $\mathbf{z} = {\mathbf{z}_1, ..., \mathbf{z}_N}$ denote a set of causal "endogenous" variables with each \mathbf{z}_i taking values in \mathbb{R} , and let $\mathbf{u} = {\mathbf{u}_1, ..., \mathbf{u}_N}$ denotes a set of mutually independent "exogenous" random variables. The latent SCM consists of a set of structural equations

$$\{\mathbf{z}_i := m_i(\mathbf{z}_{\operatorname{pa}(i)}), \mathbf{u}_i\}_{i=1}^N,\tag{B.1}$$

where $\mathbf{z}_{pa(i)}$ are the causal parents of \mathbf{z}_i and m_i are the deterministic functions that are termed "causal mechanisms". We indicate with $P_{\mathbf{u}}$ the joint distribution of the exogenous random variables, which, due to the independence hypothesis, is the product of the probability measures of the individual variables. The associated causal diagram \mathcal{G} is a directed graph with vertices \mathbf{z} and edges $\mathbf{z}_i \to \mathbf{z}_j$ iff. $\mathbf{z}_i \in \mathbf{z}_{pa(j)}$; we assume the graph \mathcal{G} to be acyclic.

The latent SCM induces a unique distribution P_z over the endogenous variables z as a pushforward of P_u via eq. (B.1). Its density p_z follows the causal Markov factorization:

$$p_{\mathbf{z}}(z) = \prod_{i=1}^{N} p_i(z_i \mid z_{\text{pa}(i)}).$$
(B.2)

Instead of directly observing the endogenous and exogenous variables z and u, we only have access to some "entangled" measurements x of z generated through a nonlinear mixing function:

Definition B.2 (Mixing function). A deterministic smooth function $f : \mathbb{R}^N \to \mathbb{R}^D$ mapping the latent vector $\mathbf{z} \in \mathbb{R}^N$ to its observable $\mathbf{x} \in \mathbb{R}^D$, where $D \ge N$ denotes the dimensionality of the observational space.

Assumption B.1 (Diffeomorphism). The mixing function f is diffeomorphic onto its image, i.e. f is C^{∞} , f is injective and $f^{-1}|_{\text{Im}(f)} : \text{Im}(f) \to \mathbb{R}^D$ is also C^{∞} .

Remark: Settings with noisy observations $(\mathbf{x} = f(\mathbf{z}) + \epsilon, \mathbf{z} \perp \epsilon)$ can be easily reduced to our denoised version by applying a standard deconvolution argument as a pre-processing step, as indicated by Lachapelle et al. (2022); Buchholz et al. (2024).

1010 C IDENTIFIABILITY THEORY

In addition to the general results for latent variable identification presented in § 3, we compare in App. C.1 different granularity of latent variable identification and show their transitions through certain assumptions on the causal model or mixing function. Afterward, App. C.2 discusses the identification level of a causal graph depending on the granularity of latent variable identification under certain structural assumptions. Proofs are deferred to App. E.

1017 C.1 ON THE GRANULARITY OF IDENTIFICATION

Different levels of identification can be achieved depending on the degree of underlying invariance and data symmetry. Below, we present three standard identifiability definitions from the CRL literature, each providing a stronger identification result than block-identifiability (Defn. 3.1).

Definition C.1 (Block affine-identifiability). Let \hat{z} be the learned representation, for a subset $A \subseteq [N]$ it satisfies that:

$$\hat{\mathbf{z}}_{\pi(A)} = D \cdot \mathbf{z}_A + \mathbf{b},\tag{C.1}$$

where $D \in \mathbb{R}^{|A| \times |A|}$ is an invertible matrix, $\pi(A)$ denotes the index permutation of A, then \mathbf{z}_A is block affine-identified by $\hat{\mathbf{z}}_{\pi(A)}$.

1034

1035

1036 1037

1039

1040

1044 1045

1048

1058

1064

1067

1068

1069

1070



Figure 2: Relations between different identification classes (Defns. 3.1 and C.1 to C.3). Some CRL works proposed a more fine-grained classification of identifiability concepts with slightly different terminology, which we omit here for readability.

Definition C.2 (Element-identifiability). The learned representation $\hat{\mathbf{z}} \in \mathbb{R}^N$ satisfies that:

$$\hat{\mathbf{z}} = \mathbf{P}_{\pi} \cdot h(\mathbf{z}),\tag{C.2}$$

where $\mathbf{P}_{\pi} \in \mathbb{R}^{N \times N}$ is a permutation matrix, $h(\mathbf{z}) := (h_1(\mathbf{z}_1), \dots h_N(\mathbf{z}_N)) \in \mathbb{R}^N$ is an elementwise diffeomorphism.

Definition C.3 (Affine-identifiability). The learned representation $\hat{\mathbf{z}} \in \mathbb{R}^N$ satisfies that:

$$\hat{\mathbf{z}} = \mathbf{\Lambda} \cdot \mathbf{P}_{\pi} \cdot \mathbf{z} + \mathbf{b},\tag{C.3}$$

where $\mathbf{P}_{\pi} \in \mathbb{R}^{N \times N}$ is a permutation matrix, $\Lambda \in \mathbb{R}^{N \times N}$ is a diagonal matrix with nonzero diagonal entries.

Remark: Block affine-identifiability (Defn. C.1) is defined by Ahuja et al. (2023), stating that the 1049 learned representation \hat{z} is related to the ground truth latents z through some sparse matrix with 1050 zero blocks. Defn. C.2 indicates element-wise identification of latent variables up to individual 1051 diffeomorphisms. Element-identifiability for the latent variable identification together with the 1052 graph identifiability (Defn. C.4) is defined as \sim_{CRL} -identifiability (von Kügelgen et al., 2024, 1053 Defn. 2.6), perfect identifiability (Varici et al., 2024a, Defn. 3). Affine identifiability (Defn. C.3) 1054 describes when the ground truth latent variables are identified up to permutation, shift, and linear 1055 scaling. In many CRL works, affine identifiability (Defn. C.3) is also termed as follows: perfect 1056 identifiability under linear transformation (Varici et al., 2024b, Defn. 1), CD-equivalence (Zhang et al., 2024a, Defn. 1), disentanglement (Lachapelle et al., 2022, Defn. 3). 1057

Proposition C.1 (Granularity of identification). *Affine-identifiability (Defn. C.3) implies element-identifiability (Defn. C.2) and block affine-identifiability (Defn. C.1) while element-identifiability and block affine-identifiability implies block-identifiability (Defn. 3.1).*

Proposition C.2 (Transition between identification levels). *The transition between different levels* of latent variable identification (Fig. 2) can be summarized as follows:

- (i) Element-level identifiability (Defns. C.2 and C.3) can be obtained from block-wise identifiability (Defns. 3.1 and C.1) when each individual latent constitutes an invariant block;
- (ii) Identifiability up to an affine transformation (Defns. C.1 and C.3) can be obtained from general identifiability on arbitrary diffeomorphism (Defns. 3.1 and C.2) by additionally assuming that both the ground truth mixing function and decoder are finite degree polynomials of the same degree.

1071 **Discussion.** We note that the granularity of identifiability results is primarily determined by the 1072 strength of invariance and parametric assumptions (such as those on mixing functions or causal models) rather than by the specific algorithmic choice. For example, for settings that can achieve element-identifiability (von Kügelgen et al., 2024), affine-identifiability results can be obtained 1074 by additionally assuming *finite degree polynomial* mixing function (proof see App. E). Similarly, 1075 one reaches element-identifiability from block-identifiability by enforcing invariance properties on each latent component (Yao et al., 2023, Thm. 3.8) instead of having only one multivariate 1077 invariant block (von Kügelgen et al., 2021). Tab. 4 provides an overview of recent identifiability 1078 results along with their corresponding invariance and parametric assumptions, illustrating the 1079 direct relationship between these assumptions and the level of identifiability they achieve.

1080 C.2 IDENTIFYING THE CAUSAL GRAPH

In addition to latent variable identification, another goal of causal representation learning is to infer the underlying latent dependency, namely the causal graph structure. Hence, we restate the standard definition of graph identifiability in causal representation learning.

Definition C.4 (Graph-identfiability). The estimated graph $\hat{\mathcal{G}}$ is isomorphic to the ground truth \mathcal{G} through a bijection $h: V(\mathcal{G}) \to V(\hat{\mathcal{G}})$ in the sense that two vertices $\mathbf{z}_i, \mathbf{z}_j \in V(\mathcal{G})$ are adjacent in \mathcal{G} if and only if $h(\mathbf{z}_i), h(\mathbf{z}_j) \in V(\hat{\mathcal{G}})$ are adjacent in $\hat{\mathcal{G}}$.

We remark that the "faithfulness" assumption (Pearl, 2009, Defn. 2.4.1) is a standard assumption in the CRL literature, commonly required for graph discovery. We restate it as follows:

Assumption C.1 (Faithfulness (or Stability)). P_z is a faithful distribution induced by the latent SCM (Defn. B.1) in the sense that P_z contains no extraneous conditional independence; in other words, the only conditional independence relations satisfied by P_z are those given by $\{\mathbf{z}_i \perp \mathbf{z}_{nd(i)} | \mathbf{z}_{pa(i)}\}$ where $\mathbf{z}_{nd(i)}$ denotes the non-descends of \mathbf{z}_i .

As indicated by Defn. C.4, the preliminary condition of identifying the causal graph is to have an element-wise correspondence between the vertices in the ground truth graph \mathcal{G} (i.e., the ground truth latents) and the vertices of the estimated graph. Therefore, the following assumes that the learned encoders G (Defn. 3.2) achieve element-identifiability (Defn. C.2), that is, for each $\mathbf{z}_i \in \mathbf{z}$, we have a differmorphism $h_i : \mathbb{R} \to \mathbb{R}$ such that $\hat{\mathbf{z}}_i = h_i(\mathbf{z}_i)$. However, additional assumptions are needed to identify the graph structure: either on the source of invariance or on the parametric form of the latent causal model.

1102 Graph identification via interventions. Under the element-identifiability (Defn. C.2) of the latent 1103 variables z, the causal graph structure \mathcal{G} can be identified up to its isomorphism (Defn. C.4), given multi-domain data from *paired perfect* interventions per-node (von Kügelgen et al., 2024; Varici 1104 et al., 2024a). Using data generated from *imperfect* interventions is generally insufficient to identify 1105 the direct edges in the causal graph. It can only identify the ancestral relations, i.e., up to the transi-1106 tive closure of \mathcal{G} (Brehmer et al., 2022; Zhang et al., 2024a). Unfortunately, even imposing the linear 1107 assumption on the latent SCM does not provide a solution (Squires et al., 2023). Nevertheless, by 1108 adding sparsity assumptions on the causal graph \mathcal{G} and polynomial assumption on the mixing func-1109 tion f, Zhang et al. (2024a) has shown isomorphic graph identifiability (Defn. C.4) under imperfect 1110 intervention per node. In general, access to the interventions is necessary for graph identification if 1111 one is uncomfortable making other parametric assumptions about the graph structure. Conveniently, 1112 in this setting, the graph identifiability is linked with that of the variables since the latter leverages 1113 the invariance induced by the intervention.

Graph identification via parametric assumptions. It is well known in causal discovery that the additive noise model (Hoyer et al., 2008) is identifiable under certain mild assumptions (Zhang & Hyvärinen, 2010; 2009). In the following, we assume an additive exogenous noise in the latent SCM (Defn. B.1):

Assumption C.2 (Additive noise). The endogenous variable $\mathbf{z}_i \in \mathbb{R}$ in the previously defined latent SCM (Defn. B.1) relates to the corresponding exogenous noise variable $\mathbf{u}_i \in \mathbb{R}$ through additivity. Namely, the causal mechanism (eq. (B.1)) can be rewritten as:

1122

$$\{\mathbf{z}_i = m_i(\mathbf{z}_{\text{pa}(i)}) + \mathbf{u}_i\}.$$
(C.4)

As a generalization of the additive noise model, the post-nonlinear acyclic causal model (Zhang & Hyvärinen, 2010, Sec. 2) allows extra nonlinearity on the top of the additive causal mechanism, providing additional flexibility on the latent model assumption:

Definition C.5 (Post-nonlinear acyclic causal model). The following causal mechanism describes a post-nonlinear acyclic causal model:

1129 1130

$$\mathbf{z}_i = h_i(m_i(\mathbf{z}_{\mathsf{pa}(i)}) + \mathbf{u}_i),\tag{C.5}$$

where $h_i : \mathbb{R} \to \mathbb{R}$ is a diffeomorphism and m_i is a non-constant function.

Assume the latent variable \mathbf{z}_i is element-wise identified through a bijective mapping $h_i : \mathbb{R} \to \mathbb{R}$ for all $i \in [N]$, define the estimated causal parents $\hat{\mathbf{z}}_{pa(i)} := \{h_j(\mathbf{z}_j) : \mathbf{z}_j \in \mathbf{z}_{pa(i)}\}$, then the latent

1134 SCM (Defn. B.1) is translated to a post-nonlinear acyclic causal model (Defn. C.5) because 1135

 $= h_i(\tilde{m}_i(\hat{\mathbf{z}}_{\mathsf{na}(i)}) + \mathbf{u}_i),$

$$\hat{\mathbf{z}}_{i} = h_{i}(\mathbf{z}_{i}) = h_{i}(m_{i}(\mathbf{z}_{pa(i)}) + \mathbf{u}_{i})$$

$$= h_{i}(m_{i}(\{h_{j}^{-1}(\hat{\mathbf{z}}_{j}) : \mathbf{z}_{j} \in \mathbf{z}_{pa(i)}\}) + \mathbf{u}_{i})$$

1138

1139

1141

 $\tilde{m}_i(\hat{\mathbf{z}}_{\mathrm{pa}(i)}) := m_i(\{h_i^{-1}(\hat{\mathbf{z}}_j) : \mathbf{z}_j \in \mathbf{z}_{\mathrm{pa}(i)}\}).$

(C.6)

1142 Thus, the underlying causal graph \mathcal{G} can be identified up to an isomorphism (Defn. C.4) following 1143 the approach given by Zhang & Hyvärinen (2009, Sec. 4)

1144 What happens if variables are identified in blocks? Consider the case where the latent variables 1145 cannot be identified up to element-wise diffeomorphism; instead, one can only obtain a coarse-1146 grained version of the variables (e.g., as a mixing of a block of variables (Defn. 3.1)). Nevertheless, 1147 certain causal links between these coarse-grained block variables are of interest. These block 1148 variables and their causal relations in between form a "macro" level of the original latent SCM, 1149 which is shown to be causally consistent under mild structural assumptions (Rubenstein et al., 2017, 1150 Thm. 11). In particular, the macro-level model can be obtained from the micro-level model through 1151 an exact transformation (Beckers & Halpern, 2019, Defn. 3.4) and thus produces the same causal effect as the original micro-level model under the same type of interventions, providing useful 1152 knowledge for downstream causal analysis. More formal connections are beyond the scope of this 1153 paper. Still, we see this concept of coarse-grained identification on both causal variables and graphs 1154 as an interesting avenue for future research. 1155

1156 D **RELATED WORKS** 1157

1158 This section reviews related causal representation learning works and frames them as specific instances of our theory (§ 3). These works were initially categorized into various causal representation 1159 learning types (multiview, multi-domain, multi-task, and temporal CRL) based on the level of invari-1160 ance in the data-generating process, leading to varying degrees of identifiability results (App. C.1). 1161 While the implementation of individual works may vary, the *methodological principle of aligning* 1162 representation with known data symmetries remains consistent, as shown in § 3. We begin with 1163 revisiting the data-generating process of each category and explain how they can be viewed as 1164 specific cases of the proposed invariance framework (§ 2). We then present individual identification 1165 algorithms from the CRL literature as particular applications of our theorems based on the 1166 implementation choices needed to satisfy the invariance and sufficiency constraints (Constraints 3.1 1167 and 3.2). A more detailed overview of the individual works is provided in Tab. 4.

1168

D.1 MULTIVIEW CAUSAL REPRESENTATION LEARNING 1169

1170 High-level overview. The multiview setting in causal representation learning (Daunhawer et al., 2023; Yao et al., 2023) considers multiple views that are *concurrently* generated by an overlapping 1171 subset of latent variables, and thus having *non-independently* distributed data. Multiview scenarios 1172 are often found in a partially observable setup. For example, multiple devices on a robot measure dif-1173 ferent modalities, jointly monitoring the environment through these real-time measurements. While 1174 each device measures a distinct subset of latent variables, these subsets probably still overlap as they 1175 are measuring the same system at the same time. In addition to partial observability, another way 1176 to obtain multiple views is to perform an "intervention/perturbation" (Locatello et al., 2020; von 1177 Kügelgen et al., 2021; Ahuja et al., 2022b; Brehmer et al., 2022) and collect both pre-action and 1178 post-action views on the same sample. This setting is often improperly termed "counterfactual"¹ 1179 in the CRL literature, and this type of data is termed "paired data". From another perspective, the 1180 paired setting can be cast in the partial observability scenario by considering the same latent before 1181 and after an action (mathematically modeled as an intervention) as two separate latent nodes in the 1182 causal graph, as shown by von Kügelgen et al. (2021, Fig. 1). Thus, both pre-action and post-action

¹Traditionally, counterfactual in causality refers to non-observable outcomes that are "counter to the 1184 fact" (Rubin, 2005). The works we refer to here represent pre- and post-actions that affect some latent vari-1185 ables but not all. This can be mathematically expressed as a counterfactual in an SCM but is conceptually 1186 different as both pre- and post-action outcomes are realized (Liu et al., 2023). The "counterfactual" terminology silently implies that this is a strong assumption, but nuance is needed and it can in fact be much weaker 1187 than an intervention.

views are partial because neither of them can observe pre-action and post-action latents simultaneously. These works assume the latents that are not affected by the action remain constant, an assumption that is relaxed in temporal CRL works. See App. D.3 for more discussion in this regard.

Data generating process. In the following, we introduce the data-generating process of a multiview setting in the flavor of the invariance principle as introduced in § 2. We consider a set of views $\{\mathbf{x}^k\}_{k\in[K]}$ with each view $\mathbf{x}^k \in \mathcal{X}^k$ generated from some latents $\mathbf{z}^k \in \mathcal{Z}^k$. Let $S_k \subseteq [N]$ be the index set of generating factors for the view \mathbf{x}^k , we define $\mathbf{z}_j^k = 0$ for all $j \in [N] \setminus S_k$ to represent the uninvolved partition of latents. Each entangled view \mathbf{x}^k is generated by a view-specific mixing function $f_k : \mathcal{Z}^k \to \mathcal{X}^k$:

$$\mathbf{x}^k = f_k(\mathbf{z}^k) \quad \forall k \in [K] \tag{D.1}$$

Define the joint overlapping index set $A := \bigcap_{k \in [K]} S_k$, and assume $A \subseteq [N]$ is a non-empty subset 1199 of [N]. Then the value of the sharing partition \mathbf{z}_A remain *invariant* for all observables $\{\mathbf{x}^k\}_{k \in [K]}$ 1200 on a *sample level*. By considering the joint intersection A, we have *one single* invariance property 1201 $\iota : \mathbb{R}^{|A|} \to \mathbb{R}^{|A|}$ in the invariance set \mathfrak{I} ; and this invariance property ι emerges as the identity map id on $\mathbb{R}^{|A|}$ in the sense that $\mathrm{id}(\mathbf{z}_A^k) = \mathrm{id}(\mathbf{z}_A^{k'})$ and thus $\mathbf{z}_A^k \sim_{\iota} \mathbf{z}_A^{k'}$ for all $k, k' \in [K]$. Note that Defn. 2.1 (ii) is satisfied because any transformation h_k that involves other components \mathbf{z}_q with 1202 1203 1204 $q \notin A$ violates the equality introduced by the identity map. For a subset of observations $V_i \subseteq [K]$ 1205 with at least two elements $|V_i| > 1$, we define the latent intersection as $A_i := \bigcap_{k \in V_i} S_k \subseteq [N]$, then 1206 for each non-empty intersection A_i , there is a corresponding invariance property $\iota_i : \mathbb{R}^{|A_i|} \to \mathbb{R}^{|A_i|}$ 1207 which is the identity map specified on the subspace $\mathbb{R}^{|A_i|}$. By considering all these subsets $\mathcal{V} :=$ 1208 $\{V_i \subseteq [K] : |V_i| > 1, |A_i| > 0\}$, we obtain a set of invariance properties $\mathfrak{I} := \{\iota_i : \mathbb{R}^{|A_i|} \to \mathbb{R}^{|A_i|}\}$ 1209 that satisfy Asm. 2.1. 1210

1211 **Identification algorithms.** Many multiview works (von Kügelgen et al., 2021; Daunhawer et al., 2023; Yao et al., 2023) employ the L_2 loss as a regularizer to enforce sample-level invariance on 1212 the invariant partition, cooperated with some sufficiency regularizer to preserve sufficient informa-1213 tion about the observables (Constraint 3.2). Aligned with our theory (Thm. 3.1), these works have 1214 shown block-identifiability on the invariant partition of the latents across different views. Follow-1215 ing the same principle, there are certain variations in the implementations to enforce the invariance 1216 principle, e.g. Locatello et al. (2020) directly average the learned representations from paired data 1217 $q(\mathbf{x}^1), q(\mathbf{x}^2)$ on the shared coordinates before forwarding them to the decoder; Ahuja et al. (2022b) 1218 enforces L_2 alignment up to a learnable sparse perturbation δ . As each latent component constitutes 1219 a single invariant block in the training data, these two works element-identifies (Defn. C.2) the latent 1220 variables, as explained by Proposition C.2.

1222 D.2 MULTI-ENVIRONMENT CAUSAL REPRESENTATION LEARNING

High-level overview. Multi-environment / interventional CRL considers data generated from mul-1223 tiple environments with respective environment-specific data distributions; hence, the considered 1224 data is *independently* but *non-identically distributed*. In the scope of causal representation learning, 1225 multi-environment data is often instantiated through interventions on the latent structured causal 1226 model (von Kügelgen et al., 2021; Zhang et al., 2024a; Buchholz et al., 2024; Squires et al., 2023; 1227 Varici et al., 2023; 2024b;a). Recently, Ahuja et al. (2024) provides a more general identifiability 1228 statement where multi-environment data is not necessarily originated from interventions; instead, 1229 they can be individual data distributions that preserve certain symmetries such as marginal invari-1230 ance or support invariance (Ahuja et al., 2024).

Data generating process The following presents the data generating process described in most interventional causal representation learning works. Formally, we consider a set of *non-identically* distributed data $\{P_{\mathbf{x}^k}\}_{k \in [K]}$ that are collected from multiple environments (indexed by $k \in [K]$) with a shared mixing function $f : \mathbf{x}^k = f(\mathbf{z}^k)$ (Defn. B.2) satisfying Asm. B.1 and a shared latent SCM (Defn. B.1). Let k = 0 denote the non-intervened environment and $\mathcal{I}_k \subseteq [N]$ denotes the set of intervened nodes in k-th environment, the latent distribution $P_{\mathbf{z}^k}$ is associated with the density

1221

1198

1239

$$p_{\mathbf{z}^k}(z^k) = \prod_{j \in \mathcal{I}_k} \tilde{p}(z_j^k \mid z_{\mathsf{pa}(j)}^k) \prod_{j \in [N] \setminus \mathcal{I}_k} p(z_j^k \mid z_{\mathsf{pa}(j)}^k), \tag{D.2}$$

where we denote by p the original density and by \tilde{p} the intervened density. Interventions naturally introduce various distributional invariances that can be utilized for latent variable identification: Under the intervention \mathcal{I}_k in the k-th environment, we observe that both (1) the marginal distribution of \mathbf{z}_A with $A := [N] \setminus \mathrm{TC}(\mathcal{I}_k)$, with TC denoting the transitive closure and (2) the score $[S(\mathbf{z}^k)]_{A'} := \nabla_{\mathbf{z}_{A'}^k} \log p_{\mathbf{z}^k}$ on the subset of latent components $A' := [N] \setminus \overline{\mathrm{pa}}(\mathcal{I}_k)$ with $\overline{\mathrm{pa}}(\mathcal{I}_k) := \{j : j \in \mathcal{I}_k \cup \mathrm{pa}(\mathcal{I}_k)\}$ remain *invariant* across the observational and the k-th interventional environment. Formally, under intervention \mathcal{I}_k , we have

• Marginal invariance:

$$p_{\mathbf{z}^0}(z_A^0) = p_{\mathbf{z}^k}(z_A^k) \qquad A := [N] \setminus \mathrm{TC}(\mathcal{I}_k); \tag{D.3}$$

• Score invariance:

1251 1252

$$[S(\mathbf{z}^0)]_{A'} = [S(\mathbf{z}^k)]_{A'} \qquad A' := [N] \setminus \overline{\operatorname{pa}}(\mathcal{I}_k). \tag{D.4}$$

According to our theory Thm. 3.1, we can block-identify both \mathbf{z}_A , \mathbf{z}'_A using these invariance principles (eqs. (D.3) and (D.4)). Since most interventional CRL works assume at least one intervention per node (Squires et al., 2023; Zhang et al., 2024a; von Kügelgen et al., 2024; Varici et al., 2024a; 2023; Buchholz et al., 2024; Ahuja et al., 2023), more fine-grained variable identification results, such as element-wise identification (Defn. C.2) or affine-identification (Defn. C.3), can be achieved by combining multiple invariances from these per-node interventions, as we elaborate below.

Identifiability with one intervention per node. By applying Thm. 3.1, we demonstrate that latent causal variables z can be identified up to element-wise diffeomorphism (Defn. C.2) under single node *imperfect* intervention per node, given the following assumption.

Assumption D.1 (Topologically ordered interventional targets). Specifying Asm. 2.1 in the interventional setting, we assume there are exactly N environments $\{k_1, \ldots, k_N\} \subseteq [K]$ where each node $j \in [N]$ undergoes one imperfect intervention in the environment $k_j \in [K]$. The interventional targets $1 \leq \cdots \leq N$ preserve the topological order, meaning that $i \leq j$ only if there is a directed path from node i to node j in the underlying causal graph \mathcal{G} .

Remark: Asm. D.1 is directly implied by Asm. 2.1 as we need to know which environments fall into 1268 the same equivalence class. We believe that identifying the topological order is another subproblem 1269 orthogonal to identifying the latent variables, which is often termed "uncoupled/non-aligned prob-1270 lem" (Varici et al., 2024a; von Kügelgen et al., 2024). As described by Zhang et al. (2024a), the 1271 topological order of unknown interventional targets can be recovered from single-node imperfect 1272 intervention by iteratively identifying the interventions that target the source nodes. This iterative 1273 identification process may require additional assumptions on the mixing functions (Zhang et al., 1274 2024a; Ahuja et al., 2023; Varici et al., 2023; 2024b; Squires et al., 2023) and the latent structured 1275 causal model (Buchholz et al., 2024; Squires et al., 2023), or on the interventions, such paired per-1276 fect interventions per node (von Kügelgen et al., 2024; Varici et al., 2024a).

Corollary D.1 (Identifiability from single node interventions per node (von Kügelgen et al., 2021)). *Given N environments* $\{k_1, \ldots, k_N\} \subseteq [K]$ *satisfying Asm. D.1, the ground truth latent variables* **z** *can be identified up to element-wise diffeomorphism (Defn. C.2) by combining both marginal and score invariances (eqs. (D.3) and (D.4)) under our framework (Thm. 3.1).*

1281 The proof for Cor. D.1 is included in App. E.5. Upon element-wise identification from single-1282 node intervention per node, existing works often provide more fine-grained identifiability results by 1283 incorporating other parametric assumptions on the mixing functions (Varici et al., 2023; Ahuja et al., 1284 2023; Zhang et al., 2024a; Squires et al., 2023). This is explained by Proposition C.2, as element-1285 wise identification can be refined to affine-identification (Defn. C.3) given additional parametric 1286 assumptions on the mixing functions. However, note that under the milder setting of *imperfect* 1287 intervention per node, the full graph is not identifiable without further assumptions. See (Zhang et al., 2024a) for more details. 1288

Identifiability with two interventions per node Current literature in interventional CRL targeting the general nonparametric setting (Varici et al., 2024a; von Kügelgen et al., 2024) typically assumed a pair of *sufficiently different* perfect interventions per node. Thus, any latent variable $z_j, j \in [N]$, as an interventional target, is uniquely shared by a pair of interventional environment $k, k' \in [K]$, forming an invariant partition $A_i = \{j\}$ constituting of individual latent node $j \in [N]$. Note that this invariance property on the interventional target induces the following distributional property:

$$S(\mathbf{z}^{k}) - S(\mathbf{z}^{k'})]_{j} \neq 0 \qquad \text{only if} \qquad \mathcal{I}_{k} = \mathcal{I}_{k'} = \{j\}.$$
(D.5)

According to Thm. 3.1, each latent variable can thus be identified separately, giving rise to elementwise identification, as shown by (Varici et al., 2024a; von Kügelgen et al., 2024).

Identifiability under multiple distributions. More recently, Ahuja et al. (2024) explains previous 1299 interventional identifiability results from a general weak distributional invariance perspective. In a 1300 nutshell, a set of variables z_A can be block-identified if certain invariant distributional properties 1301 hold: The invariant partition z_A can be block-identified (Defn. 3.1) from the rest by utilizing the 1302 marginal distributional invariance or invariance on the support, mean or variance. Ahuja et al. 1303 (2024) additionally assume the mixing function to be finite degree polynomial, which leads to block-1304 affine identification (Defn. C.1), whereas we can also consider a general nonparametric setting; they 1305 consider one single invariance set, which is a special case of Thm. 3.1 with one joint ι -property. 1306

1307 **Identification algorithms.** Instead of iteratively enforcing the invariance constraint across the majority of environments as described in Cor. D.1, most single-node interventional works develop 1308 equivalent constraints between pairs of environments to optimize. For example, the marginal 1309 invariance (eq. (D.3)) implies the marginal of the source node is changed only if it is intervened 1310 upon, which is utilized by Zhang et al. (2024a) to identify latent variables and the ancestral relations 1311 simultaneously. In practice, Zhang et al. (2024a) propose a regularized loss that includes Maximum 1312 Mean Discrepancy(MMD) between the reconstructed "counterfactual" data distribution and the 1313 interventional distribution, enforcing the distributional discrepancy that reveals graphical structure 1314 (e.g., detecting the source node). Similarly, by enforcing sparsity on the score change matrix, 1315 Varici et al. (2023) restricts only score changes from the intervened node and its parents. In the 1316 nonparametric case, von Kügelgen et al. (2024) optimize for the invariant (aligned) interventional 1317 targets through model selection, whereas Varici et al. (2024a) directly solve the constrained 1318 optimization problem formulated using score differences. Considering a more general setup, Ahuja 1319 et al. (2024) provides various invariance-based regularizers as plug-and-play components for any losses that enforce a sufficient representation (Constraint 3.2). 1320

1321

1322 D.3 TEMPORAL CAUSAL REPRESENTATION LEARNING

High-level overview. Temporal CRL (Lippe et al., 2022a; 2023; 2022b; Yao et al., 2022a;b; 1323 Lachapelle et al., 2022; 2024; Li et al., 2024a;b) focuses on retrieving latent causal structures from 1324 time series data, where the latent causal structured is typically modeled as a Dynamic Bayesian 1325 Network (DBN) (Dean & Kanazawa, 1989; Murphy, 2002). Existing temporal CRL literature has 1326 developed identifiability results under varying sets of assumptions. A common overarching assump-1327 tion is to require the Dynamic Bayesian Network to be first-order Markovian, allowing only causal 1328 links from t - 1 to t, eliminating longer dependencies (Lippe et al., 2022b; 2023; 2022a; Yao et al., 1329 2022b). While many works assume that there is no instantaneous effect, restricting the latent com-1330 ponents of \mathbf{z}^t to be mutually dependent (Lippe et al., 2022b; Yao et al., 2022b; Lippe et al., 2023), 1331 some approaches have lifted this assumption and prove identifiability allowing for instantaneous 1332 links among the latent components at the same timestep (Lippe et al. (2022a)). 1333

Data generating process. We present the data generating process followed by most temporal causal 1334 representation works and explain the underlying latent invariance and data symmetries. Let $\mathbf{z}^t \in \mathbb{R}^N$ 1335 denotes the latent vector at time t and $\mathbf{x}^t = f(\mathbf{z}^t) \in \mathbb{R}^D$ the corresponding entangled observable 1336 with $f : \mathbb{R}^N \to \mathbb{R}^D$ the shared mixing function (Defn. B.2) satisfying Asm. B.1. The actions \mathbf{a}^t with 1337 cardinality $|\mathbf{a}^t| = N$ mostly only target a subset of latent variables while keeping the rest untouched, 1338 following its default dynamics (Lippe et al., 2022b; 2023; Lachapelle et al., 2022; 2024). Intuitively, 1339 these actions \mathbf{a}^t can be interpreted as a component-wise indicator for each latent variable $\mathbf{z}_i^t, j \in [N]$ 1340 stating whether \mathbf{z}_j follows the default dynamics $p(\mathbf{z}_j^{t+1} \mid \mathbf{z}^t)$ or the modified dynamics induced by 1341 the action \mathbf{a}_{i}^{t} . From this perspective, the non-intervened causal variables at time t can be considered 1342 the invariant partition under our formulation, denoted by $\mathbf{z}_{A_t}^t$ with the index set A_t defined as $A_t :=$ 1343 $\{j : \mathbf{a}_j = 0\}$. Note that this invariance can be considered as a generalization of the multiview case 1344 because the realizations z_j^t, z_j^{t+1} are not exactly identical (as in the multiview case) but are related 1345 via a default transition mechanism $p(\mathbf{z}_j^{t+1} \mid \mathbf{z}^t)$. To formalize this intuition, we define $\tilde{\mathbf{z}}^t := \mathbf{z}^t \mid \mathbf{a}^t$ 1346 as the conditional random vector conditioning on the action \mathbf{a}^t at time t. For the non-intervened 1347 1348 partition $A_t \subseteq [N]$ that follows the default dynamics, the transition model should be invariant: 1349

$$p(\mathbf{z}_{A_t}^t \mid \mathbf{z}^{t-1}) = p(\tilde{\mathbf{z}}_{A_t}^t \mid \mathbf{z}^{t-1}), \tag{D.6}$$

1350 which gives rise to a non-trivial distributional invariance property (Defn. 2.1). Note that the invari-1351 ance partition A_t could vary across different time steps, providing a set of invariance properties 1352 $\mathfrak{I} := \{\iota_t : \mathbb{R}^{|A_t|} \to \mathcal{M}_t\}_{t=1}^T$, indexed by time t. Given by Thm. 3.1, all invariant partitions $\mathbf{z}_{A_t}^t$ can be block-identified; furthermore, as shown in Proposition 3.3, the complementary variant 1353 1354 partition can also be identified under an invertible encoder and mutual independence within \mathbf{z}^t 1355 (here conditioning on the previous time step z^{t-1}), aligning with the identification results without 1356 instantaneous effect, i.e. there is no causal link between variables at the same time step (Lippe et al., 2022b; Yao et al., 2022b; Lachapelle et al., 2022; 2024). On the other hand, temporal causal 1357 variables with instantaneous effects are shown to be identifiable only if "instantaneous parents" (i.e., 1358 nodes affecting other nodes instantaneously) are cut by actions (Lippe et al., 2022a), reducing to 1359 the setting without instantaneous effect where the latent components at t are mutually independent. 1360 Upon invariance, more fine-grained latent variable identification results, such as element-wise 1361 identifiability, can be obtained by incorporating additional technical assumptions, such as the sparse 1362 mechanism shift (Lachapelle et al., 2022; 2024; Li et al., 2024b) and parametric latent causal 1363 model (Yao et al., 2022b; Klindt et al., 2021; Khemakhem et al., 2020). 1364

Identification algorithms. From a high level, the distributional invariance (eq. (D.6)) indicates 1365 full explainability and predictability of $\mathbf{z}_{A_t}^t$ from its previous time step \mathbf{z}^{t-1} , regardless of the action at. In principle, this invariance principle can be enforced by directly maximizing the 1367 information content of the proposed default transition density between the learned representation 1368 $p(\hat{\mathbf{z}}_{A_{+}}^{t} | \hat{\mathbf{z}}^{t-1})$ (Lippe et al., 2022a;b). In practice, the invariance regularization is often incorporated 1369 together with the predictability of the variant partition conditioning on actions, implemented 1370 as a KL divergence between the observational posterior $q(\hat{\mathbf{z}}^t \mid \mathbf{x}^t)$ and the transitional prior 1371 $p(\hat{\mathbf{z}}^t | \hat{\mathbf{z}}^{t-1}, \mathbf{a}^t)$ (Lachapelle et al., 2022; 2024; Klindt et al., 2021; Yao et al., 2022a; Lippe et al., 1372 2023), estimated using variational Bayes (Kingma & Welling, 2013) or normalizing flow (Rezende 1373 & Mohamed, 2015).

1374

1375 D.4 MULTI-TASK CAUSAL REPRESENTATION LEARNING

1376 High-level overview. Multi-task causal representation learning aims to identify latent causal 1377 variables via external supervision, in this case, the label information of the same instance for 1378 various tasks. Previously, multi-task learning (Caruana, 1997; Zhang & Yang, 2018) has been mostly studied outside the scope of identifiability, mainly focusing on domain adaptation and 1379 out-of-distribution generalization. One of the popular ideas that was extensively used in the context 1380 of multi-task learning is to leverage interactions between different tasks to construct a generalist 1381 model that is capable of solving all classification tasks and potentially better generalizes to unseen 1382 tasks (Zhu et al., 2022; Bai et al., 2022). Recently, Lachapelle et al. (2023); Fumero et al. (2024) systematically studied under which conditions the latent variables can be identified in the multi-task 1384 scenario and correspondingly provided identification algorithms. 1385

Data generating process. The multi-task causal representation learning considers a *supervised* setup: Given a latent SCM as defined in Defn. B.1, we generate the observable $\mathbf{x} \in \mathbb{R}^D$ through some mixing function $f : \mathbb{R}^N \to \mathbb{R}^D$ satisfying Asm. B.1. Given a set of task $\mathcal{T} = \{t_1, \ldots, t_k\}$, and let $\mathbf{y}^k \in \mathcal{Y}_k$ denote the corresponding task label respect to the task t_k . Each task only *directly* depends on a subset of latent variables $S_k \subseteq [N]$, in the sense that the label \mathbf{y}^k can be expressed as a function that contains all and only information about the latent variable \mathbf{z}_{S_k} :

1392

$$\mathbf{y}^k = r_k(\mathbf{z}_{S_k}),\tag{D.7}$$

1393 1394 where $r : \mathbb{R}^{|S_k|} \to \mathcal{Y}_k$ is some deterministic function which maps the latent subspace $\mathbb{R}^{|S_k|}$ to the 1395 task-specific label space \mathcal{Y}_k , which is often assumed to be linear and implemented using a linear 1396 readout in practice (Lachapelle et al., 2023; Fumero et al., 2024). For each task $t_k, k \in [K]$, we 1397 observe the associated data distribution $P_{\mathbf{x},\mathbf{y}^k}$. Consider two different tasks $t_k, t_{k'}$ with $k, k' \in [K]$, 1398 the corresponding data \mathbf{x}, \mathbf{y}^k and $\mathbf{x}, \mathbf{y}^{k'}$ are *invariant* in the intersection of task-related features \mathbf{z}_A 1399 with $A = S_k \cap S_{k'}$. Formally, let $r_k^{-1}(\{\mathbf{y}^k\})$ denotes the pre-image of \mathbf{y}^k , for which it holds

1400
1400
1401
$$r_k^{-1}(\{\mathbf{y}^k\})_A = r_{k'}^{-1}(\{\mathbf{y}^{k'}\})_A,$$
(D.8)

showing alignment on the shared partition of the task-related latents. In the ideal case, each latent component $j \in [N]$ is *uniquely shared* by a subset of tasks, all factors of variation can be fully disentangled, which aligns with the theoretical claims by Lachapelle et al. (2023); Fumero et al. (2024). 1404 **Identification algorithms.** We remark that the *sharing* mechanism in the context of multi-task 1405 learning fundamentally differs from that of multiview setup, thus resulting in different learning 1406 algorithms. Regarding learning, the shared partition of task-related latents is enforced to align 1407 up to the linear equivalence class (given a linear readout) instead of sample level L_2 alignment. 1408 Intuitively, this invariance principle can be interpreted as a soft version of the that in the multiview case. In practice, under the constraint of perfect classification, one employs (1) a sparsity constraint 1409 on the linear readout weights to enforce the encoder to allocate the correct task-specific latents 1410 and (2) an information-sharing term to encourage reusing latents across various tasks. Equilibrium 1411 can be obtained between these two terms only when the shared task-specific latent is element-wise 1412 identified (Defn. C.2). Thus, this soft invariance principle is jointly implemented by the sparsity 1413 constraint and information sharing regularization (Fumero et al., 2024, Sec. 2.1). 1414

1415 D.5 DOMAIN GENERALIZATION

1416 High-level overview. Domain generalization aims at out-of-distribution performance. That is, learn-1417 ing an optimal encoder and predictor that performs well at some unseen test domain that preserves 1418 the same data symmetries as in the training data. At a high level, domain generalization represen-1419 tation learning (Sagawa et al., 2019; Zhang et al., 2017; Ganin et al., 2016; Arjovsky et al., 2020; Krueger et al., 2021) considers a similar framework as introduced for interventional CRL, with inde-1420 pendent but non-identically distributed data, but additionally incorporated with external supervision 1421 and focusing more on model robustness perspective. While interventional CRL aims to identify the 1422 true latent factors of variations (up to some transformation), domain generalization learning focuses 1423 directly on *out-of-distribution* prediction, relying on some invariance properties preserved under the 1424 distributional shifts. Due to the non-causal objective, new methodologies are motivated and tested 1425 on real-world benchmarks (e.g., VLCS (Fang et al., 2013), PACS (Li et al., 2017), Office-Home 1426 (Venkateswara et al., 2017), Terra Incognita (Beery et al., 2018), DomainNet (Peng et al., 2019)) 1427 and could inspire future real-world applicability of causal representation learning approaches. 1428

Data generating process. The problem of domain generalizations is an *extension of supervised* 1429 *learning* where training data from multiple environments are available (Blanchard et al., 2011). An 1430 environment is a dataset of i.i.d. observations from a joint distribution $P_{\mathbf{x}^k, \mathbf{y}^k}$ of the observables 1431 $\mathbf{x}^k \in \mathbb{R}^D$ and the label $\mathbf{y}^k \in \mathbb{R}$. The label $\mathbf{y}^k \in \mathbb{R}^m$ only depends on the invariant latents through 1432 a linear regression structural equation model (Ahuja et al., 2022a, Assmp. 1), described as follows: 1433

У

- 1434
- 1435 1436

1443

1444

$$\mathbf{y}^{k} = \mathbf{w}^{*} \mathbf{z}_{A}^{k} + \epsilon_{k}, \, \mathbf{z}_{A}^{k} \perp \epsilon_{k}$$
$$\mathbf{x}^{k} = f(\mathbf{z}^{k})$$
(D.9)

where $\mathbf{w}^* \in \mathbb{R}^{D \times m}$ represents the ground truth relationship between the label \mathbf{y}^k and the invariant latents \mathbf{z}_A^k . ϵ_k is some white noise with bounded variance and $f : \mathbb{R}^N \to \mathbb{R}^D$ denotes the 1437 1438 shared mixing function for all $k \in [K]$ satisfying Asm. B.1. The set of environment distributions 1439 $\{P_{\mathbf{x}^k,\mathbf{v}^k}\}_{k\in[K]}$ generally differ from each other because of interventions or other distributional 1440 shifts such as covariates shift and concept shift. However, as the relationship between the invariant 1441 latents and the labels \mathbf{w}^* and the mixing mechanism f are shared across different environments, the 1442 optimal risk remains invariant in the sense that

$$\mathcal{R}_k^*(\mathbf{w}^* \circ f^{-1}) = \mathcal{R}_{k'}^*(\mathbf{w}^* \circ f^{-1}), \tag{D.10}$$

1445 where \mathbf{w}^* denotes the ground truth relation between the invariant latents \mathbf{z}_A^k and the labels \mathbf{y}^k and 1446 f^{-1} is the inverse of the diffeomorphism mixing f (see eq. (D.9)). Note that this is a non-trivial ι 1447 property as the labels y^k only depend on the invariant latents z_A^k , thus satisfying Defn. 2.1 (ii). 1448

Identification algorithms. Different distributional invariance are enforced by interpolating and 1449 extrapolating across various environments. Among the countless contribution to the literature, 1450 mixup (Zhang et al., 2017) linearly interpolates observations from different environments as a robust 1451 data augmentation procedure, Domain-Adversarial Neural Networks (Ganin et al., 2016) support 1452 the main learning task discouraging learning domain-discriminant features, Distributionally Robust 1453 Optimization (DRO) (Sagawa et al., 2019) replaces the vanilla Empirical Risk objective minimizing 1454 only with respect to the worst modeled environment, Invariant Risk Minimization (Arjovsky et al., 2020) combines the Empirical Risk objective with an invariance constraint on the gradient, and 1455 Variance Risk Extrapolation (Krueger et al., 2021, V-REx), similar in spirit combines the empirical 1456 risk objective with an invariance constraint using the variance among environments. For a more 1457 comprehensive review of domain generalization algorithms, see Zhou et al. (2022).

1458 D.6 FURTHER EXPLANATIONS FOR TAB. 4

1476 1477

General clarification. Tab. 4 summarizes special cases of our invariance framework. For each work, we present their technical assumptions, the type of invariance, the implementation for the invariance and the sufficiency regularizers (to satisfy Constraints 3.1 and 3.2), and the type of identifiability they achieve. Note that this table is by no means exhaustive. Also, we omit some additional results and technical assumptions of individual papers for readability. A list of paragraphs is provided below for further clarification, as referenced in Tab. 4.

(a) Single-node intervention and parametric assumptions. Many existing CRL works that consider single node intervention per node require additional parametric assumptions, either on the mixing function (Varici et al., 2023; Zhang et al., 2024a) or the latent causal model (Buchholz et al., 2024) or both (Squires et al., 2023), thus achieving (at least) element-wise identifiability (Defn. C.2). We conjecture these additional parametric assumptions serve two purposes: (1) to identify valid topological order of the interventional targets, as required by Asm. D.1 for Cor. D.1 (2) to get a more fine-grained identification level of affine transformation, as explained by Proposition C.2.

1472 In the following, we restate the definition of linear latent SCM for reference:

Definition D.1 (Linear latent SCM (Squires et al., 2023; Buchholz et al., 2024)). The latent variables
 z follows a linear SCM with Gaussian noise in the sense that

$$\mathbf{z} = A\mathbf{z} + \Gamma^{1/2}\epsilon,\tag{D.11}$$

1478 where Γ is a diagonal matrix with positive entries, A encodes the underlying causal graph G and the 1479 ϵ is the standard Gaussian noise. For the sake of simplicity, we often define $B := \Gamma^{-1/2}(\text{Id} - A)$ 1480 such that $\mathbf{z} = B^{-1}\epsilon$ to explicitly map from the exogenous noise ϵ to the latent variables \mathbf{z} . We use 1481 B_k to denote this matrix for the domain k.

1482 (b) Multi-node intervention and linear mixing. Recently, Varici et al. (2024b) extends previous 1483 interventional CRL works to unknown multi-node interventions and achieves identifiability under the assumption of a linearly independent intervention signature matrix $M_{\text{int}} \in \{0,1\}^{N \times K}$ 1484 with each column k represents the intervened node in this environment k. The row-wise linear 1485 independence of M_{int} implies that each latent variable must have been intervened at least once. Let $M \in \{0,1\}^{N \times N}$ represent a submatrix of M_{int} with *linearly independent* columns. By 1486 1487 multiplying M with its adjoint transpose $\operatorname{adj}^{\mathsf{T}}(M)$, one obtains a matrix where each column has 1488 only one non-zero component. Applying the same transformation to the score change, this problem 1489 is reduced to a similar setting as a single node intervention per node, which can be intuitively 1490 explained using the same distributional invariance principle introduced earlier (App. D.2). 1491

(c) Paired single-node intervention per node under nonparametric assumptions. In the nonparametric settings, several works (von Kügelgen et al., 2024; Varici et al., 2024a) have shown element-wise latent variable identification under sufficiently different paired perfect intervention per node. By having two sufficiently different interventions per node, one introduces invariance on the interventional target across these paired interventional environments. This invariance property can be enforced using the score differences (Varici et al., 2024a) or algorithmically by performing model selection (von Kügelgen et al., 2024), as elaborated in App. D.2.

(d) Variant latents identification under independence. While some papers states main identification results on the variant partition, it can be explained by Thm. 3.1 and Proposition 3.3 stating that the variant block can be identified under independence and invertible encoder. For example, Wendong et al. (2024, Thm. 4.5) shows block-identifiability on the intervened (variant) latents under (Wendong et al., 2024, Assumption 4.4) of block-wise independence between the invariant and variant blocks.

(e) Invariance regularizers in multitask CRL Under the assumption of knowing the number of latent variables, Lachapelle et al. (2023) solves a bi-level optimization problem, enforcing $L_{2,1}$ sparsity on individual task readouts in the inner problem. Coupled with a backbone shared across all tasks, this implicitly encourages discovering the ground truth overlapping partition of task support. Fumero et al. (2024) lifted the constraint of assuming the known number of latents by incorporating an additional information-sharing regularizer, as explained in (Fumero et al., 2024, Sec. 2.1).

(f) Invariance regularizers in domain generalization. While Sagawa et al. (2019) directly optimize for the worst-case risk, a link can be drawn between this objective and the risk invariance:

1512 Given a pair of linear head w and encoder g shared across [K] domains, let the order of risks be 1513 $\mathcal{R}^{\pi_1} \geq \mathcal{R}^{\pi_2} \dots \mathcal{R}^{\pi_K}$. Since \mathcal{R}^{π_1} is lower bounded by \mathcal{R}^{π_2} the minimum of the training objective 1514 in Sagawa et al. (2019) (max_{k \in [K]} \mathcal{R}^k(w, g)) is obtained when $\mathcal{R}^{\pi_1} = \mathcal{R}^{\pi_2}$. Then we have $\mathcal{R}^{\pi_1} =$ 1515 $\mathcal{R}^{\pi_2} \geq \dots \geq \mathcal{R}^{\pi_K}$, and the next minimum will be obtained when $\mathcal{R}^{\pi_1} = \mathcal{R}^{\pi_2} = \mathcal{R}^{\pi_3}$, and so on so 1516 forth. The optimization procedure stops when the risks are equally minimized across all domains.

(Krueger et al., 2021) minimizes variance between domain risks to enforce the risk invariance. We formally show these two are equivalent in the following. Note that the invariance principle for risk alignment can be formulated as

1520

1524

1525 1526

$$(\mathcal{R}_k - \mathcal{R}_{k'})^2 \tag{D.12}$$

According to Zhang et al. (2012), variance can be equivalently expressed as pair-wise distances between the samples. Hence, we can reformulate the risk variance term in (Sagawa et al., 2019) as follows:

$$\operatorname{Var}\left[\mathcal{R}\right] = \frac{1}{K^2} \sum_{k,k' \in [K]} \frac{1}{2} \left(\mathcal{R}_k - \mathcal{R}_{k'}\right)^2,$$

showing that the variance regularization in (Krueger et al., 2021) enforces risk invariance.

1528 D.7 NOTABLE CASES NOT DIRECTLY COVERED BY THE THEORY

1529 Some works not listed in Tab. 4 cannot yet be directly explained by our invariance frameworks but 1530 are rather loosely connected. One representative line of work (Lachapelle et al., 2022; Zheng et al., 1531 2022; Xu et al., 2024; Lachapelle et al., 2024) relies on the sparsity assumption in the latent depen-1532 dency to achieve latent variable and graph identification. This assumption is closely related to the 1533 sparse mechanism shift hypothesis in causal representation learning (Schölkopf et al., 2021), stating 1534 small distributional changes should not affect all causal variables but only a small subset of these. 1535 Note that the sparsity constraint is often formulated as the estimator (either for the graph (Lachapelle et al., 2023; 2024) or of the latents (Xu et al., 2024)) should be at least sparse as the ground truth one, 1536 maximizing the cardinality of the unaffected (invariant) part. Some theoretical results do not rely on 1537 multiple data pockets that share certain invariance properties but directly employ specific properties 1538 within the observational data, such as independent support (Ahuja et al., 2023), or shared cluster 1539 membership (Khemakhem et al., 2020; Kivva et al., 2022). Some works (Zhang et al., 2024b) fol-1540 low an orthogonal proof technique originating from the nonlinear ICA with auxiliary variable line 1541 of work (Hyvarinen et al., 2019). Their proofs often rely on linear independence derived from the 1542 statistical diversity of various underlying data distributions instead of shared invariance properties. 1543 Our framework thus does not trivially include them.

1545 E PROOFS

This section includes formal proofs for the theoretical statements of the paper.

1548 E.1 ASSUMPTION JUSTIFICATION

We justify the Defn. 2.1 (ii) by showing negative results under violation of this assumption, i.e., trivially invariant latent variables are not identifiable.

Proposition E.1 (General non-identifiability of trivially invariant latent variables). Consider the setup in Thm. 3.1, w.l.o.g we assume $\Im = \{\iota\}$ and ι is trivial in the sense that assumption (ii) in Defn. 2.1 is violated. Then, the corresponding invariant partition \mathbf{z}_A^k is not identifiable for any $k \in [K]$.

¹⁵⁵⁶ *Proof.* We provide a counter example as follows: Define a trivial ι -property as "if the first component is greater than zero on $A = \{1\}$ of some two dimensional latents z". Formally,

1558 1559

1544

$$\iota(\mathbf{z}_1) = \mathbf{1}[\mathbf{z}_1 > 0].$$

1560 Consider a mixing function f = id and an invertible encoder $g(\mathbf{x}) = g(f(\mathbf{z})) = [\mathbf{z}_1 + \mathbf{z}_2, \mathbf{z}_2]$ 1561 satisfying the sufficiency constraint (Constraint 3.2). Define $h_1 = h_2 = [g \circ f]_A$. Then for some 1562 realizations z, \tilde{z} with $z_1 + z_2 > 0$ and $\tilde{z}_1 + \tilde{z}_2 > 0$ we have $\iota(h(\mathbf{z})) = \iota(h(\tilde{\mathbf{z}}))$. However, h_1, h_2 can 1563 not disentangle \mathbf{z}_1 , showing non-identifiability for the invariant partition \mathbf{z}_A .

- 1564
- **1565** Link between Defn. 2.1 (ii) and interventional discrepancy. In the following, we elaborate how Defn. 2.1 (ii) resembles the most common assumption in interventional causal representation

Ċ

learning, the interventional discrepancy (Wendong et al., 2024; Varici et al., 2024a). Note that this assumption may termed differently as *sufficient variability* (von Kügelgen et al., 2024; Lippe et al., 2022b), *interventional regularity* (Varici et al., 2023; 2024b), but the mathematical formulation remain the same. We begin with restating this assumption:

1570 1571 1572 1573 Assumption E.1 (Interventional discrepancy (Wendong et al., 2024)). Given $k \in [K]$, let p_{t_k} denote the causal mechanism of the intervened variable \mathbf{z}_{t_k} with $t_k \in [N]$. We say a stochastic intervention \tilde{p}_k satisfies interventional discrepancy if

$$\frac{\partial \log p_{t_k}}{\partial \mathbf{z}_{t_k}} (\mathbf{z}_{t_k} \mid \mathbf{z}_{\text{pa}(t_k)}) \neq \frac{\partial \log \tilde{p}_{t_k}}{\partial \mathbf{z}_{t_k}} (\mathbf{z}_{t_k} \mid \mathbf{z}_{\text{pa}(t_k)}) \quad \text{almost everywhere } (a.e.).$$

1575 1576

1574

1577 *Proof.* We show that any cases violating the interventional discrepancy assumption also vio-1578 lates Defn. 2.1 (ii) and vice versa. Suppose for a contradiction that there exists $t_k \in [N]$ that is 1579 intervened in environment $k \in [K]$, and there is a non-empty interior $U \subset \mathbb{R}$ with non-zero measure 1580 where the interventional discrepancy is violated, i.e., for all $z_{t_k} \in U$, it holds

$$\frac{\partial \log p_{t_k}}{\partial z_{t_k}}(\mathbf{z}_{t_k} \mid \mathbf{z}_{\operatorname{pa}(t_k)}) = \frac{\partial \log \tilde{p}_{t_k}}{\partial z_{t_k}}(\mathbf{z}_{t_k} \mid \mathbf{z}_{\operatorname{pa}(t_k)})$$
(E.1)

1583

1587

1590

1591

1594

1596

1598

1610

1611

1612

1613

1614

1615

1581

Under a single node imperfect intervention, the complementary set of the transitive closure of t_k , i.e., $A := [N] \setminus TC(t_k)$ remain marginally invariant:

$$\iota(\mathbf{z}_A) = p_{\mathbf{z}_A} = \tilde{p}_{\mathbf{z}_A}$$

1589 W.l.o.g, we assume $A = \{1, \dots, t_k - 1\}$, define a function $h : \mathbb{R}^N \to \mathbb{R}^{|A|}$ with

 $h(\mathbf{z}) = [\mathbf{z}_1, \dots, \mathbf{z}_{t_k-2}, \mathbf{z}_{t_k}]$

that omits the t_k -1-th component of z but includes the variant component t_k . Note that the marginal of z_{t_k} after intervention remains invariant within U because

$$p(\mathbf{z}_{t_k}) = \int p_{t_k}(\mathbf{z}_{t_k} | \mathbf{z}_{pa(t_k)}) p(\mathbf{z}_{pa(t_k)}) d\mathbf{z}_{pa(t_k)} \qquad pa(t_k) \in A$$
$$= \int p_{t_k}(\mathbf{z}_{t_k} | \mathbf{z}_{pa(t_k)}) \tilde{p}(\mathbf{z}_{pa(t_k)}) d\mathbf{z}_{pa(t_k)} \qquad eq. \text{ (E.1) and both } p_k, \tilde{p}_k \text{ pdfs}$$
$$= \int \tilde{p}_{t_k}(\mathbf{z}_{t_k} | \mathbf{z}_{pa(t_k)}) \tilde{p}(\mathbf{z}_{pa(t_k)}) d\mathbf{z}_{pa(t_k)}$$
$$= \tilde{p}(\mathbf{z}_{t_k}).$$

Therefore, we have $\iota(h(\mathbf{z})) = \iota(h(\tilde{\mathbf{z}}))$ (with $\tilde{\mathbf{z}}$ noting the latent vectors under intervention) contradicting Defn. 2.1 (ii). The other direction (violating Defn. 2.1 (ii) implies violating Asm. E.1) can be proved using the same example.

1607 1608 E.2 PROOF FOR THM. 3.1

1609 Our proof consists of the following steps:

- 1. We construct the optimal encoders G^* (Defn. 3.2) and selectors Φ^* (Defn. 3.4) that solves the constrained optimization problem in Thm. 3.1.
- 2. We show that, for any invariance property $\iota_i \in \mathfrak{I}$ and any observation \mathbf{x}^k in the corresponding ι_i -equivalent subset \mathbf{x}_{V_i} , the selected representation $\phi^{(i,k)} \oslash g_k(\mathbf{x}^k)$ cannot contain any other information than the invariant partition $\mathbf{z}_{A_i}^k$.
- 1616 1617 1618 1619 3. Lastly, we prove that selected representation $\phi^{(i,k)} \oslash g_k(\mathbf{x}^k)$ relates to the ground truth invariant partition $\mathbf{z}_{A_i}^k$ through a diffeomorphism $h_k : \mathbb{R}^{|A_i|} \to \mathbb{R}^{|A_i|}$ for all invariance property $\iota_i \in \mathfrak{I}$ and for any observable \mathbf{x}^k from the ι_i -equivalent subset \mathbf{x}_{V_i} ; in other words, $\phi^{(i,k)} \oslash g_k(\mathbf{x}^k)$ block-identifies $\mathbf{z}_{A_i}^k$ in the sense of Defn. 3.1.

Lemma E.1 (Existence of optimal encoders and selectors). Consider a set of observables $S_{\mathbf{x}} = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^K\} \in \mathcal{X}$ generated from § 2 satisfying Asm. 2.1, then there exists optimal encoders G^* (Defn. 3.2) and selectors Φ^* (Defn. 3.4) which satisfy both Constraints 3.1 and 3.2.

1624 *Proof.* The optimal encoders can be constructed as the set of the inverse of the ground truth mixing functions: 1626 $C^* = \{f^{-1}\}$

$$G^* = \{f_k^{-1}\}_{k \in [K]},\tag{E.2}$$

 f_{k}^{-1} is smooth and invertible following Asm. B.1. By definition, for each $k \in [K]$, we have:

$$f_k^{-1}(\mathbf{x}^k) = \mathbf{z}^k \in \mathcal{Z}^k.$$
(E.3)

1631 Next, we define the optimal selector $\Phi^* = \{\phi^{(i,k)}\}_{i \in [n_{\mathfrak{I}}], k \in [K]}$ such that for all $i \in n_{\mathfrak{I}}, k \in [K]$, it holds

$$\phi^{(i,k)} \oslash \mathbf{z}^k = \mathbf{z}_{A_i}^k. \tag{E.4}$$

Thus, the invariance constraint (Constraint 3.1) is trivially satisfied as given by § 2. The optimal encoder f_k^{-1} is smooth and invertible following Asm. B.1 so the sufficiency constraint (Constraint 3.2) is also satisfied. Hence, we have shown the optimum of the constrained optimization problem in Thm. 3.1 exists.

1639 1640 1641 1641 1642 1643 1644 Lemma E.2 (Invariant component isolation). Consider the same set of observables $S_{\mathbf{x}}$ as introduced in Lemma E.1, then for any set of smooth encoders G (Defn. 3.2), Φ (Defn. 3.4) that satisfy the invariance condition (Constraint 3.1), the learned representation $\phi^{(i,k)} \oslash g_k(\mathbf{x}^k)$ can only be dependent on the invariant latent variables $\mathbf{z}_{A_i}^k := {\mathbf{z}_j^k : j \in A_i}$, not any non-invariant variables \mathbf{z}_q^k with $q \in A_i^c := [N] \setminus A_i$.

Proof. This proof directly follows Defn. 2.1 (ii). Define

$$h_k^i := \phi^{(i,k)} \oslash g_k \circ f_k \quad k \in [K].$$
(E.5)

1648 1649 By Constraint 3.1, for all $\iota_i \in \mathfrak{I}$, we have

$$\iota_i(h_k^i(\mathbf{z}^k)) = \iota_i(h_{k'}^i(\mathbf{z}^{k'})) \quad a.s. \quad \forall k \neq k' \in [K].$$
(E.6)

According to Defn. 2.1 (ii), for all $i \in [n_{\mathfrak{I}}], k \in V_i, h_k^i$ cannot directly depends on any other latent component \mathbf{z}_q with $q \notin A_i$. Therefore, we have shown that h_k^i is a function of $\mathbf{z}_{A_i}^k$, for all $i \in [n_{\mathfrak{I}}], k \in V_i$.

Theorem 3.1 (Identifiability of multiple invariant blocks). Consider a set of observables $S_{\mathbf{x}} = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^K\} \in \mathcal{X}$ generated from § 2 satisfying Asm. 2.1. Let G, Φ be the set of smooth encoders (Defn. 3.2) and selectors (Defn. 3.4) that satisfy Constraints 3.1 and 3.2, then the invariant component $\mathbf{z}_{A_i}^k$ is block-identified (Defn. 3.1) by $\phi^{(i,k)} \oslash g_k$ for all $\iota_i \in \mathfrak{I}, k \in [K]$.

1660

1627

1629 1630

1633 1634

1645

1646 1647

1650 1651

1661 *Proof.* Lem. E.1 verifies that there exists such optimum which satisfies both invariance and suffi-1662 ciency conditions (Constraints 3.1 and 3.2). Following Lem. E.2, the composition $\phi^{(i,k)} \otimes g_k$ can 1663 only encode information related to the invariant latent subset A_i specified by the invariance property 1664 $\iota_i \in \mathfrak{I}$ for all $k \in V_i$. As given by Constraint 3.2, $\phi^{(i,k)} \otimes g_k$ contain all information the ground 1665 truth invariant latents \mathbf{z}_{A_i} for i with $k \in V_i$. Therefore, the selected representation $\phi^{(i,k)} \otimes g_k(\mathbf{x}^k)$ 1666 relates to the ground truth invariant partition \mathbf{z}_{A_i} through some diffeomorphism, i.e., \mathbf{z}_{A_i} is blocked-1667 identified by $\phi^{(i,k)} \otimes g_k(\mathbf{x}^k)$ for all invariance property $\iota_i \in \mathfrak{I}$ and observable $k \in V_i$.

1668

1669 E.3 PROOFS FOR GENERALIZATION OF VARIANT LATENTS

Proposition 3.2 (General non-identifiability of variant latent variables). Consider the setup in Thm. 3.1, let $A := \bigcup_{i \in [n_{\mathfrak{I}}]} A_i$ denote the union of block-identified latent indices and $A^c := [N] \setminus A$ the complementary set where no ι -invariance $\iota \in \mathfrak{I}$ applies, then the variant latents \mathbf{z}_{A^c} cannot be identified. 1674 1675 *Proof.* We provide a simple counter example with two latent variables $\mathbf{z} = [\mathbf{z}_1, \mathbf{z}_2]$, with the mixing 1676 function f being the identity map id. W.l.o.g. we assume the invariant partition to be $A = \{1\}$. 1676 According to Thm. 3.1, the invariant latent variable can be identified up to a certain bijection h: 1677 $\mathbb{R} \to \mathbb{R}$. Let $\hat{\mathbf{z}}$ be the estimated representation:

$$\hat{\mathbf{z}} = [h(\mathbf{z}_1), \mathbf{z}_2 - \mathbf{z}_1] \tag{E.7}$$

1680 with the estimated mixing function $\hat{f} : \mathbb{R}^2 \to \mathbb{R}^2$:

$$\hat{f}(\hat{\mathbf{z}}) = [h^{-1}(\hat{\mathbf{z}}_1), \hat{\mathbf{z}}_2 + h^{-1}(\hat{\mathbf{z}}_1)],$$
 (E.8)

then we obtain the same observations $\hat{f}(\hat{z}) = f(z)$ whereas \hat{z}_2 consists of a mixing of z_1 and z_2 , showing the variant latent variable z_2 can not be identified.

Proposition 3.3 (Identifiability of variant latent under independence). Consider an optimal encoder $g \in G^*$ and optimal selector $\phi \in \Phi^*$ from Thm. 3.1 that jointly identify an invariant block \mathbf{z}_A (we omit subscriptions k, i for simplicity), then $\mathbf{z}_{A^c}(A^c := [N] \setminus A)$ can be identified by the complementary encoding partition $(1 - \phi) \oslash g$ only if

(i) g is invertible in the sense that $I(\mathbf{x}, g(\mathbf{x})) = H(\mathbf{x})$;

(*ii*) \mathbf{z}_{A^c} *is independent on* \mathbf{z}_A .

Proof. We start by showing the sufficiency of conditions (i) and (ii). The mutual information between the observation $\mathbf{x} \in S_{\mathbf{x}}$ and the optimal encoder $g \in G^*$ from Thm. 3.1 writes:

$$I(\mathbf{x}, g(\mathbf{x})) = H(\mathbf{x}) - H(\mathbf{x} \mid g(\mathbf{x})),$$

following condition (i) in Proposition 3.3, the second term (conditional entropy) must equal zero: $H(\mathbf{x} \mid g(\mathbf{x})) = 0.$

1700 Writing the $\mathbf{x} = f(\mathbf{z}_A, \mathbf{z}_{A^c})$, we have

$$H(\mathbf{x} \mid g(\mathbf{x})) = H(f(\mathbf{z}_A, \mathbf{z}_{A^c}) \mid g(\mathbf{x})) = H(\mathbf{z}_A, \mathbf{z}_{A^c} \mid g(\mathbf{x})),$$

because the mixing function f is deterministic as given by Defn. B.2.

Note that $g(\mathbf{x})$ can be decomposed into two separate partitions: $\phi \oslash g(\mathbf{x}), (1 - \phi) \oslash g(\mathbf{x})$; thus we can write the conditional entropy as

1706 1707

1708

1709 1710

1713 1714

1717

1720

1701

1678 1679

1681 1682

1693

$$\begin{aligned} H(\mathbf{x} \mid g(\mathbf{x})) &= H(\mathbf{z}_A, \mathbf{z}_{A^c} \mid \phi \oslash g(\mathbf{x}), (1 - \phi) \oslash g(\mathbf{x})) \\ &= H(\mathbf{z}_{A^c} \mid \mathbf{z}_A, \phi \oslash g(\mathbf{x}), (1 - \phi) \oslash g(\mathbf{x})) + H(\mathbf{z}_A \mid \phi \oslash g(\mathbf{x}), (1 - \phi) \oslash g(\mathbf{x})) \end{aligned}$$

Given that $\phi \oslash g(\mathbf{x})$ block identifies \mathbf{z}_A , $(1 - \phi) \oslash g(\mathbf{x})$) cannot contain any information about \mathbf{z}_A , hence we can simplify the second term as

Using the additional mutual independence assumption between \mathbf{z}_A and \mathbf{z}_{A^c} (Proposition 3.3 (ii)), we can rewrite the first term as

$$H(\mathbf{z}_{A^{c}} \mid (1-\phi) \oslash g(\mathbf{x})).$$

 $H(\mathbf{z}_A \mid \phi \oslash g(\mathbf{x}))$

1718 As a result, the condition entropy $H(\mathbf{x} \mid g(\mathbf{x}))$ can be decomposed as

$$H(\mathbf{x} \mid g(\mathbf{x})) = H(\mathbf{z}_A \mid \phi \oslash g(\mathbf{x})) + H(\mathbf{z}_{A^c} \mid (1 - \phi) \oslash g(\mathbf{x})) = 0.$$

1721 Since $H(\mathbf{z}_A \mid \phi \oslash g(\mathbf{x})) = 0$ following Constraint 3.2, the second term also must be zero, i.e., 1722 $H(\mathbf{z}_{A^c} \mid (1-\phi) \oslash g(\mathbf{x})) = 0$, which is satisfied only if $(1-\phi) \oslash g(\mathbf{x})$ is a invertible function of 1723 \mathbf{z}_{A^c} . That is, $(1-\phi) \oslash g(\mathbf{x})$ block-identifies \mathbf{z}_{A^c} .

Next, we show both (i) and (ii) are necessary conditions for the statement in Proposition 3.3 using two counterexamples.

For (i). We consider the following scenario where condition (ii) holds, but condition (i) is violated. Given a non-invertible encoder $g \in G^*$ where $g(\mathbf{x})_A = h(\mathbf{z}_A)$ for some diffeomorphism h (implied by Thm. 3.1) and $g(\mathbf{x})_{A^c} = \vec{0} \in \{0\}^{|A^c|}$. The encoder g is non-invertible because the A^c partition does not contain any information about \mathbf{z}_{A^c} but only constant zeros. This can be formally verified by the mutual information between the observable \mathbf{x} and its encoding $g(\mathbf{x})$:

$$I(\mathbf{x}, g(\mathbf{x})) = I(f(\mathbf{z}_A, \mathbf{z}_{A^c}), [h(\mathbf{z}_A), \vec{0}]) = H(\mathbf{z}_A)$$

which is smaller than $H(\mathbf{x})$ because the ground truth partition \mathbf{z}_{A^c} contains independent information of \mathbf{z}_A . Formally, the mapping between the ground truth \mathbf{z}_{A^c} and the representations $g(\mathbf{x})_{A^c}$ becomes a constant function

1732

$$h_0: \mathcal{Z}_{A^c} \to \{0\}^{|A^c|}, \ h_0(\mathbf{z}_{A^c}) = \vec{0}$$
 (E.9)

which is clearly not a diffeomorphism, indicating $g(\mathbf{x})_{A^c}$ does not block-identify \mathbf{z}_{A^c} . Thus, condition (i) is shown to be necessary for the statement in Proposition 3.3.

For (ii). Now, assume the complementary partition \mathbf{z}_{A^c} depends on the identified partition \mathbf{z}_A , thereby violating condition (ii). For example, let $\mathbf{z}_A = \mathbf{z}_{A^c}$ for some observable $\mathbf{x} = f(\mathbf{z}_A, \mathbf{z}_{A^c}) = f(\mathbf{z}_A)$ (Note that here we slightly abuse the notation f for simplification). Consider the same encoder $g \in G^*$ as described in the counterexample for (i), i.e., $g(\mathbf{x})_A = h(\mathbf{z}_A)$ for some diffeomorphism h (implied by Thm. 3.1) and $g(\mathbf{x})_{A^c} = \vec{0} \in \{0\}^{|A^c|}$. However, note that this encoder g is invertible because

 $I(\mathbf{x}, g(\mathbf{x})) = I(f(\mathbf{z}_A), [h(\mathbf{z}_A), \vec{0}]) = H(\mathbf{z}_A) = H(\mathbf{x}).$

1748 1749 Nevertheless, the mapping between the encoding $g(\mathbf{x})_{A^c}$ and the ground truth latents \mathbf{z}_{A^c} remains the constant zero mapping h_0 (eq. (E.9)), which fails to block-identify \mathbf{z}_{A^c} .

To this end, we have shown both condition (i) and (ii) and necessary and sufficient conditions for the block-identifiability of \mathbf{z}_{A^c} which completes the proof.

1753 1754

1755

1747

E.4 PROOFS FOR GRANULARITY OF LATENT VARIABLE IDENTIFICATION

Proposition C.1 (Granularity of identification). Affine-identifiability (Defn. C.3) implies element-identifiability (Defn. C.2) and block affine-identifiability (Defn. C.1) while element-identifiability and block affine-identifiability implies block-identifiability (Defn. 3.1).

1759

Proof. The diagonal matrix Λ in eq. (C.3) is invertible and thus also a diffeomorphism h (eq. (C.2)). Hence, affine-identifiability implies element-identifiability. Affine-identifiability provides identification results with block-size one thus implies block affine-identifiability. On the other hand, block affine-identifiability is block-identifiability with affine bijection h and element-identifiability defines a special case of block-identifiability where each latent component \mathbf{z}_i is an individual block.

1765 1766

1767

1768

1769

1770 1771

1772

1773

Proposition C.2 (Transition between identification levels). *The transition between different levels of latent variable identification (Fig. 2) can be summarized as follows:*

- (i) Element-level identifiability (Defns. C.2 and C.3) can be obtained from block-wise identifiability (Defns. 3.1 and C.1) when each individual latent constitutes an invariant block;
- (ii) Identifiability up to an affine transformation (Defns. C.1 and C.3) can be obtained from general identifiability on arbitrary diffeomorphism (Defns. 3.1 and C.2) by additionally assuming that both the ground truth mixing function and decoder are finite degree polynomials of the same degree.

1774 1775

1781

1776 1777 *Proof.* The proof for (i) is trivial in the sense that identification of block with size one boils down to the identification on the element level. (ii) directly follows Ahuja et al. (2023, Thm. 4.4) and Zhang et al. (2024a, Lem. 1), stating that when both ground truth mixing function and decoder are finite degree polynomials of the same degree, the *invertible* encoder learns a representation that is affine linear to the ground truth latents, i.e., $\hat{z} = L \cdot z + b$ with $L \in \mathbb{R}^{N \times N}$.



1792Figure 3: Causal Model for generic partially annotated scientific experiment: T treatment, W1793experimental settings, X high-dimensional observation, Y outcome, S annotation flag. Figure and caption adapted from (Cadei et al., 2024, Fig. 1)1795Figure 3: Causal Model for generic partially annotated scientific experiment: T treatment, W1794experimental settings, X high-dimensional observation, Y outcome, S annotation flag. Figure and caption adapted from (Cadei et al., 2024, Fig. 1)



(a) Grooming (blue to focal)

(b) No Action

Figure 4: Examples of high-dimensional observations X with corresponding annotated social behaviour Y (grooming). Figure and caption adapted from (Cadei et al., 2024, Fig. 2)

1799 E.5 PROOF FOR COR. D.1

Corollary D.1 (Identifiability from single node interventions per node (von Kügelgen et al., 2021)). *Given N environments* $\{k_1, \ldots, k_N\} \subseteq [K]$ *satisfying Asm. D.1, the ground truth latent variables* **z** *can be identified up to element-wise diffeomorphism (Defn. C.2) by combining both marginal and score invariances (eqs.* (D.3) *and* (D.4)) *under our framework (Thm. 3.1).*

180

1798

1805 *Proof.* We consider a coarse-grained version of the underlying causal graph consisting of a block-1806 node $\mathbf{z}_{[N-1]} := {\mathbf{z}_1, \dots, \mathbf{z}_{N-1}}$ and the leaf node \mathbf{z}_N with $\mathbf{z}_{[N-1]}$ causing \mathbf{z}_N (i.e., $\mathbf{z}_{[N-1]} \to \mathbf{z}_N$). 1807 We first select a pair of environments $V = \{0, k_N\}$ consisting of the observational environment and 1808 the environment where the leaf node \mathbf{z}_N is intervened upon. According to eq. (D.3), the marginal 1809 *invariance* holds for the partition A = [N-1], implying identification on $\mathbf{z}_{[N-1]}$ from Thm. 3.1. 1810 At the same time, when considering the set of environments $V' = \{0, k_1, \dots, k_{N-1}\}$, the leaf node N is the only component that satisfy *score* invariance across all environments V', because N is not 1811 the parent of any intervened node (also see (Varici et al., 2023, Lemma 4)). So here we have another 1812 invariant partition $A' = \{N\}$, implying identification on \mathbf{z}_N (Thm. 3.1). By jointly enforcing the 1813 marginal and score invariance on A and A' under a sufficient encoder (Constraint 3.2), we identify 1814 both $\mathbf{z}_{[N-1]}$ as a block and \mathbf{z}_N as a single element. Formally, for the parental block $\mathbf{z}_{[N-1]}$, we have: 1815

$$\hat{\mathbf{z}}_{[N-1]}^k = g_{:N-1}(\mathbf{x}^k) \qquad \forall k \in \{0, k_1, \dots, k_N\}$$
(E.10)

where $g_{:N-1}(\mathbf{x}^k) := [g(\mathbf{x}^k)]_{:N-1}$ relates to the ground truth $\mathbf{z}_{[N-1]}$ through some diffeomorphism $h_{[N-1]} : \mathbb{R}^{N-1} \to \mathbb{R}^{N-1}$ (Defn. 3.1). Now, we can remove the leaf node N as follows: For each environment $k \in \{0, k_1, \dots, k_{N-1}\}$, we compute the pushforward of $P_{\mathbf{x}^k}$ using the learned encoder $g_{:N-1} : \mathcal{X}^k \to \mathbb{R}^{N-1}$:

$$P_{\hat{\mathbf{z}}_{[N-1]}^k} = g_{\#}(P_{\mathbf{x}^k})$$

Note that the estimated representations $P_{\hat{\mathbf{z}}_{[N-1]}^k}$ can be seen as a new observed data distribution for each environment k that is generated from the subgraph \mathcal{G}_{-N} without the leaf node N. Using an iterative argument, we can identify all latent variables element-wise (Defn. C.2).

1827

1823

1816 1817

¹⁸²⁸ F IMPLEMENTATION DETAILS

This section provides further details about the experiment settings of § 5, including a formal introduction to the ISTAnt dataset, highlighted open challenges (App. F.1), and additional training settings for reproducibility (App. F.2).

F.1 CASE STUDY: ISTANT

Problem. Despite the majority of causal representation learning algorithms being designed to enforce the identifiability of some latent factors and tested on controlled synthetic benchmarks, there

1836	Model/Usinen nenemeters	Value(s)
1837	Would Hyper-parameters	value(s)
1838	Encoder	DINOv2 (Oquab et al., 2023)
1839	Encoder (token)	class
1840	MLP (head): hidden layers	1
1841	MLP (head): hidden nodes	256
10/10	MLP (head): activation function	ReLU + Sigmoid output
1042	Tass	or
1843	Dropout	No
1844	Regularization	No
1845	Loss	BCELoss (with positive weighting)
1846	Loss: Positive Weight	$\sum_{i=1}^{n_s} 1 - Y_i$
1847	Loomin - Dete	$\sum_{i=1}^{n_s} Y_i$
1848	Learning Rate	0.0005
1849	Dutah Si	Adam $(p_1 = 0.9, p_2 = 0.9, \epsilon = 10^{-4})$
1850	Batch Size	128
1051	Epocns	15
1001	Seeds	range(20)

Table 2: Model and training details for the case study on ISTAnt (§ 5.1). Table adapted from (Cadei et al., 2024, Tab. 4)

1856 1857

1855

are a plethora of real-world applications across scientific disciplines requiring representation learn-1859 ing to answer causal questions (Robins et al., 2000; Samet et al., 2000; Van Nes et al., 2015; Runge, 2023). Recently, Cadei et al. (2024) introduced ISTAnt, the first real-world representation learning 1861 benchmark with a real causal downstream task (treatment effect estimation). This benchmark high-1862 lights different challenges (sources of biases) that could arise from machine learning pipelines even 1863 in the simplest possible setting of a randomized controlled trial. Videos of ants triplets are recorded, 1864 and a per-frame representation has to be extracted for supervised behavior classification to estimate the Average Treatment Effect of an intervention (exposure to a chemical substance). Beyond de-1865 sirable identification result on the latent factors (implying that the causal variables are recovered 1866 without bias), no clear algorithm has been proposed yet on minimizing the Treatment Effect Bias 1867 (TEB) (Cadei et al., 2024). One of the challenges highlighted by Cadei et al. (2024) is that in prac-1868 tice, there is both covariate and concept shifts due to the effect modification from training on a 1869 non-random subset of the RCT because, for example, ecologists do not label individual frames but 1870 whole video recordings. Figs. 3 and 4 shows the underlying causal graph and example input. 1871

Solution. Relying on our framework, we can explicitly aim for low TEB by leveraging *known data symmetries* from the experimental protocol. In fact, the causal mechanism $(P(Y^e|do(X^e = x)))$ stays invariant among the different experiment settings (i.e., individual videos or position of the petri dish). This condition can be easily enforced by existing domain generalization algorithms. For exemplary purposes, we choose Variance Risk Extrapolation (Krueger et al., 2021, V-REx), which directly enforces both the invariance sufficiency constraints (Constraints 3.1 and 3.2) by minimizing the Empirical Risk together with the risk variance inter-environments.

1879

Implementation details All training settings follow the best-performing settings from (Cadei et al., 2024), which we restate in Tab. 2 for reference.

Discussion. Interestingly, Gulrajani & Lopez-Paz (2020) empirically demonstrated that no domain generalization algorithm consistently outperforms Empirical Risk Minimization in *out-ofdistribution* prediction. However, in this application, our goal is not to achieve high out-ofdistribution accuracy but rather to identify a representation that is invariant to the effect modifiers introduced by the data labeling process. This experiment serves as a clear example of the paradigm shift of CRL via the invariance principle. While existing CRL approaches design algorithms based on specific assumptions that are often challenging to align with real-world applications, our approach begins from the application perspective. It allows for the specification of known data symmetries and desired properties of the learned representation, followed by selecting an appropriate implemen-

Parameter	Value
Mixing function	3-layer MLP
Encoder	3-layer MLP
Decoder	3-layer MLP
Hidden dim	128
Activation	Leaky-ReLU
Optimizer	Adam
Adam: learning rate	1e-4
Adam: beta1	0.9
Adam: beta2	0.999
Adam: epsilon	1e-8
Batch size	4000
Sample size	200,000
# Epochs	500

Table 3: Training setup for synthetic ablations in § 5.2.

1905 1906

1913 1914 1915

1916 1917

1925 1926 1927

1890

1892

1894 1895

1897

tation for the distance function (potentially from existing methods). Ultimately, identifiability hinges
 on the guarantee of asymptotic consistency in the estimates.

1910 F.2 Synthetic Ablation with "Ninterventions"

1911 The numerical data is generated using a linear Gaussian additive noise model as follows:

$$p(\mathbf{z}_{1}) = \mathcal{N}(\mu_{1}, \sigma_{1}^{2})$$

$$p(\mathbf{z}_{2} \mid \mathbf{z}_{1}) = \mathcal{N}(\alpha_{1} \cdot \mathbf{z}_{1} + \beta_{1}, \sigma_{2}^{2})$$

$$p(\mathbf{z}_{3} \mid \mathbf{z}_{2}) = \mathcal{N}(\alpha_{2} \cdot \mathbf{z}_{2} + \beta_{2}, \sigma_{3}^{2})$$

$$\tilde{p}(\mathbf{z}_{2}) = \mathcal{N}(\tilde{\mu}_{2}, \tilde{\sigma}_{2}^{2})$$
(F.1)

We choose $\mu_1 = 10.5, \sigma_1 = 0.8, \alpha_1 = 0.02, \beta_1 = 0, \sigma_2 = 0.5, \alpha_2 = 1, \beta_2 = 3, \sigma_3 = 1, \tilde{\sigma}_2 = 0.02$. We sample three independent $\tilde{\mu}_2$ according to a uniform distribution Unif[2,5] to validate the consistency of the identification results.

For the training, we employ a simple auto-encoder architecture implementing both encoder and decoder as 3-Layer MLP. We enforce the marginal invariance using the Max Mean Discrepancy loss (MMD) on the first and last component \hat{z}_1 , \hat{z}_3 . Formally, the objective function writes

$$\mathcal{L}(g,\hat{f}) = \mathbb{E}_{\mathbf{x},\tilde{\mathbf{x}}} \left[\left\| \hat{f}(g(\mathbf{x})) - \mathbf{x} \right\|_{2}^{2} + \left\| \hat{f}(g(\tilde{\mathbf{x}})) - \mathbf{x} \right\|_{2}^{2} \right] + \mathrm{MMD}(g(\mathbf{x})_{[1,3]}, g(\tilde{\mathbf{x}})_{[1,3]}),$$

1928 where $\mathbf{x}, \tilde{\mathbf{x}}$ denote the observational and ninterventional data, respectively.

Further training details are summarized in Tab. 3

1931 G FURTHER DISCUSSIONS AND CONNECTIONS TO OTHER FIELDS

In this paper, we take a closer look at the wide range of causal representation learning methods. 1933 Interestingly, we find that the differences between them may often be more related to "semantics" 1934 than to fundamental methodological distinctions. We identified two components involved in identi-1935 fiability results: preserving information of the data and a set of known invariances. Our results have 1936 two immediate implications. First, they provide new insights into the "causal representation learning problem," particularly clarifying the role of causal assumptions. We have shown that while learning 1938 the graph requires traditional causal assumptions such as additive noise models or access to inter-1939 ventions, identifying the causal variables may not. This is an important result, as access to causal variables is standalone useful for downstream tasks, e.g., for training robust downstream predictors or even extracting pre-treatment covariates for treatment effect estimation (Yao et al., 2024), even 1941 without knowledge of the full causal graph. Second, we have exemplified how causal representation 1942 can lead to successful applications in practice. We moved the goal post from a characterization of 1943 specific assumptions that lead to identifiability, which often do not align with real-world data, to a

general recipe that allow practitioners to specify known invariances in their problem and learn representations that align with them. In the domain generalization literature, it has been widely observed that invariant training methods often do not consistently outperform empirical risk minimization (ERM). In our experiments, instead, we have demonstrated that the specific invariance enforced by V-REx (Krueger et al., 2021) entails good performance in our causal downstream task (§ 5.1). Our paper leaves out certain settings concerning identifiability that may be interesting for future work, such as discrete variables and finite samples guarantees.

1951 One question the reader may ask, then, is "so what is exactly causal in causal representation learn-1952 ing?". We have shown that the identifiability results in typical causal representation learning are 1953 primarily based on invariance assumptions, which do not necessarily pertain to causality. We hope 1954 this insight will broaden the applicability of these methods. At the same time, we used causality as 1955 a language describing the "parameterization" of the system in terms of latent causal variables with 1956 associated known symmetries. Defining the symmetries at the level of these causal variables gives the identified representation a causal meaning, important when incorporating a graph discovery step 1957 or some other causal downstream task like treatment effect estimation. Ultimately, our representa-1958 tions and latent causal models can be "true" in the sense of (Peters et al., 2014) when they allow 1959 us to predict "causal effects that one observes in practice". Overall, our view also aligns with "phenomenological" accounts of causality (Janzing & Mejia, 2024), that define causal variables from a 1961 set of elementary interventions. In our setting too, the identified latent variables or blocks thereof 1962 are directly defined by the invariances at hand. From the methodological perspective, all is needed 1963 to learn causal variables is for the symmetries defined over the causal latent variables to entail some 1964 statistical footprint across pockets of data. If variables are available, learning the graph has a rich 1965 literature (Peters et al., 2017), with assumptions that are often compatible with learning the variables 1966 themselves. Our general characterization of the variable learning problem opens new frontiers for 1967 research in representation learning:

1968 1969

G.1 REPRESENTATIONAL ALIGNMENT AND PLATONIC REPRESENTATION

Several works (Li et al. (2015); Moschella et al. (2022); Kornblith et al. (2019); Huh et al. (2024)) 1970 have highlighted the emergence of similar representations in neural models trained independently. 1971 In Huh et al. (2024) is hypothesized that neural networks, trained with different objectives on various 1972 data and modalities, are converging toward a *shared* statistical model of reality within their represen-1973 tation spaces. To support this hypothesis, they measure the alignment of representations proposing 1974 to use a mutual nearest-neighbor metric, which measures the mean intersection of the k-nearest 1975 neighbor sets induced by two kernels defined on the two spaces, normalized by k. This metric can 1976 be an instance to the distance function in our formulation in Thm. 3.1. Despite not being optimized directly, several models in multiple settings (different objectives, data and modalities) seem to be 1978 aligned, hinting at the fact that their individual training objectives may be respecting some unknwon 1979 symmetries. A precise formalization of the latent causal model and identifiability in the context of foundational models remains open and will be objective for future research.

1981 1982 G.2 Environment Discovery

Domain generalization methods generalize to distributions potentially far away from the training, distribution, via learning representations invariant across distinct environments. However this can 1984 be costly as it requires to have label information informing on the partition of the data into environments. Automatic environment discovery (Creager et al. (2021); Arefin et al. (2024); Pezeshki et al. 1986 (2024)) attempts to solve this problem by learning to recover the environment partition. This is an in-1987 teresting new frontier for causal representation learning, discovering data symmetries as opposed to 1988 only enforcing them. For example, this would correspond to having access to multiple interventional 1989 distributions but without knowing which samples belong to the same interventional or observational 1990 distribution. Discovering that a data set is a mixture of distributions, each being a different intervention on the same causal model, could help increase applicability of causal representations to large 1992 obeservational data sets. We expect this to be particularly relevant to downstream tasks were biases 1993 to certain experimental settings are undesirable, as in our case study on treatment effect estimation from high-dimensional recordings of a randomized controlled trial.

- G.3 GEOMETRIC DEEP LEARNING
- 1997 Geometric deep learning (GDL) (Bronstein et al. (2017; 2021)) is a well estabilished learning paradigm which involves encoding a geometric understanding of data as an inductive bias in deep

learning models, in order to obtain more robust models and improve performance. One fundamental direction for these priors is to encode symmetries and invariances to different types of transfor-mations of the input data, e.g. rotations or group actions (Cohen & Welling (2016); Cohen et al. (2018)), in representational space. Our work can be fundamentally related with this direction, with the difference that we don't aim to model *explicitly* the transformations of the input space, but the invariances defined at the latent level. While an initial connection has been developed for disentan-glement Fumero et al. (2021); Higgins et al. (2018), a precise connection between GDL and causal representation learning remains a open direction. We expect this to benefit the two communities in both directions: (i) by injecting geometric priors in order to craft better CRL algorithms and (ii) by incorporating causality into successful GDL frameworks, which have been fundamentally advancing challenging real-world problems, such as protein folding (Jumper et al. (2021)).

2053	f th
2054	y c
2055	ivit
2056	ect
2057	inj
2058	me
2059	inse
2060	s ac
2061	ork
2062	M
2063	ted
2064	lis
2065	the
2066	of 1
2067	Ţ
2068	⊲,
2069	ing
2070	E
2071	Lea
2072	[u
2073	atio
2074	nta
2075	ese
2076	epr
2077	Å
2078	sal
2079	Jau
2080	L L
2081	fo
2082	ılts
2083	est
2084	V L
	ilit
	ìab
	ltif
	der
	tin
	exis.
	ofe
	Ň
	nai
	m
	ns
	Ve

ē mixing function and causal sufficiency (Markovianity) for the causal latent variables. Many listed papers depend on further technical assumptions and could yield additional results. For clarity, these are omitted; see references for details. In the table, "not assigned" means that the practical method did not directly enforce the invariance principle but considered other algorithmic designs that still implicitly preserve the data symmetries. Table 4: A non-exhaustiv

k	Causal Model	Mixing Function	Invariance	Source of invari- ance, Inv. subset A	Invariance reg.	Sufficiency reg.	Identifiability	Expl.
res al. 3, 1	linear	linear	distributional	l perfect interven- tion per node	$ rank(H^{T} \Delta_{k}H) \stackrel{!}{=} 1 $ for source nodes; linear encoder $g(\mathbf{x}) = H\mathbf{x}$, where $\Delta_{k} := B_{k}^{T} B_{k} - B_{0}^{T} B_{0}, \mathbf{z} =$ $B_{k}^{-1} \epsilon$	g invertible by assump- tion	affine-id. and partial order preserving graph-id.	(a)
ja al. 4,	nonparam.	finite- deg. poly.	marginal	single-node imper- fect interventions on variant latents	$\sum_{k, k'} \sum_{j \in A} MMD(p_{[g}^k(\mathbf{x})]_j, p_{[g}^{k'}(\mathbf{x})]_j)$	$\sum_{k} \mathbb{E}_{\mathbf{x}} \left\ \hat{f}(g(\mathbf{x}^{k})) - \mathbf{x}^{k} \right\ _{2}^{2}$	block affine- id.	1
ja al. 4,	nonparam.	finite- deg. poly.	marginal	multi-node imper- fect interventions on variant latents	$\sum_{k, k'} \sum_{j \in A} MMD(p_{[g}^k(\mathbf{x})]_j, p_{[g}^{k'}(\mathbf{x})]_j)$	$\sum_{k} \mathbb{E}_{\mathbf{x}} k \left\ \hat{f}(g(\mathbf{x}^{k})) - \mathbf{x}^{k} \right\ _{2}^{2}$	block affine- id.	 1
ja 4, al.	nonparam.	finite- deg. poly.	marginal support	imperfect interven- tions on variant la- tents	$\begin{split} \sum_{k,k'}\sum_{j\inA} \\ &\left\ \operatorname{bnd}(\hat{\mathcal{Z}}_{j}^{k}) - \operatorname{bnd}(\hat{\mathcal{Z}}_{j}^{k'}) \right\ _{2}^{2} \end{split}$	$\sum_{k} \mathbb{E}_{\mathbf{x}} k \left\ \hat{f}(g(\mathbf{x}^{k})) - \mathbf{x}^{k} \right\ _{2}^{2}$	block affine- id.	.
al. 4)	linear Gaussian	nonparam.	marginal	perfect interven- tion per node	$\frac{-\mathbb{E}_l \sim \mathcal{U}(\{0,k\}) \mathbb{E}_{\mathbf{x}^l}}{\ln\left(e^{1_{l} = k \cdot g \cdot \mathbf{k}^l}\right)}$	$\mathbb{E}_{l} \sim \mathcal{U}(\{0,k\}) \mathbb{E}_{\mathbf{x}^{l}}$ $\ln \left(e^{g_{k}}(\mathbf{x}^{l}) + 1 \right)$	affine id. + graph id.	(a)

2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114

Vork	Causal Model	Mixing Function	Invariance Source of invari- ance, Inv. subset A	Invariance reg.	Sufficiency reg.	Identifiability	Expl
'arici t al.		linoor	distributional perfect interven-	$\begin{aligned} \left\ \Delta_{\mathbf{x}}^{s}(U^{T}) \right\ _{0}, \text{For all} \\ j, k \in [N], \text{its element} \\ \left[\Delta_{\mathbf{x}}^{s}(U^{T}) \right]_{j,k} = \end{aligned}$	g invertible by assump-	affine-id. +	
2023, hm. 16)	nouparam.		tion per node	$\begin{aligned} 1([U^{T}S(\mathbf{x}^0)]_j \stackrel{P_{\mathbf{x}^0,k}}{\neq} [U^{T}S(\mathbf{x}^k)]_j \\ g(\mathbf{x}) := U^+ \mathbf{x} \end{aligned}$	tion),	graph-id.	a
/arici t al.		-	····· , imperfect interven-	$\begin{split} \ \Delta_{s}^{s}(U^{\intercal})\ _{0}, & \text{For all} \\ j, k \in [N], & \text{its element} \\ [\Delta_{s}^{s}(U^{\intercal})]_{j,k} = \end{split}$	<i>a</i> invertible by assump-	block affine-	
2023, 'hm. 13)	nonparam.	uncar	distributional tion per node	$1([U^{T}S(\mathbf{x}^{0})]_{j} \stackrel{P_{\mathbf{x}^{0},k}}{\neq} [U^{T}S(\mathbf{x}^{k})]_{j}$ $g(\mathbf{x}):=U^{+}\mathbf{x}$	tion .	id. + grapn- id.	(a)
/arici				$\min \ \Delta^{s}(g)\ _{0} \text{s.t. it is}$			
t al. 2024a, 'hm. 3)	nonparam.	nonparam.	interventional partect inter- target vention per node	diagonal. $\Delta^{s}(g)_{j,k} = \mathbb{E}\left[[S(g(\mathbf{x}^{k})) - S(g(\mathbf{x}^{k'}))]_{j} \right]$	g invertible by assump- tion	element-1a. + graph-id.	(c)
				Linear encoder $g(\mathbf{x}) = H\mathbf{x}$,			
/arici			linearly indepen-	$H_i^*\!\in\!\mathrm{im}(\Delta s_{\mathbf{x}}\mathbf{w}_i)\backslash\mathrm{span}(H_{[i-1]}^*)$			
t al. 2024b.	nonparam.	linear	distributional dent multi-node perfect interven-	such that the \dim of	g invertible by assump- tion	affine id. + graph id.	(q)
hm. 1)			tion	$\operatorname{proj}_{\operatorname{null}} \left(H^*_{[i-1]} \right) \operatorname{im}(\Delta S_{\mathbf{x}} \mathbf{w}_i)$			
				equals one.			

ork	Causal Model	Mixing Function	Invariance	Source of invari- ance, Inv. subset A	Invariance reg.	Sufficiency reg.	Identifiability	Expl.
					Linear encoder $g(\mathbf{x}) = H\mathbf{x}$,			
rici				linearly indepen-	$H_i^*\!\in\!\mathrm{im}(\Delta s_{\mathbf{x}}\mathbf{w}_i)\backslash\mathrm{span}(H_{[i-1]}^*)$		hlock affine.	
al. 124h	nonparam.	linear	distribution	al dent multinode al immerfect interven-	such that the \dim of	g invertible by assump-	id. + graph	(q)
n. 2)				tion	$\operatorname{proj}_{\operatorname{null}} \left(H^*_{[i-1]} \right) \operatorname{im}(\Delta S_{\mathbf{x}} \mathbf{w}_i)$		id.	
					equals one.			
ang al. 24a)	nonparam.	finite- deg. poly.	distribution	al imperfect interven- tion per node	$-\sum_{k} MMD(q_{\mathbf{x}^{k}}, p_{\mathbf{x}^{k}})$ where \mathbf{x}^{k} the generated "counterfactual" pair through VAE	$-\sum_k \mathbb{E}_{\mathbf{x}^k} \log p(\mathbf{x}^k g(\mathbf{x}^k))$	affine-id. + graph id.	(a)
ndong al. 24, n. 4.5)	nonparam.	nonparam.	marginal	marginal invari- ance from multiple fat-hand interven- tions on the same set of interven- tional targets I , invariant partition $A:=[N]\setminus I$	model selection	$-\sum_k \log p_{oldsymbol{ heta}}^k(\mathbf{x}^k)$	block-id. (known graph)	(p)
gel- et al. 24, n. 4.1)	nonparam.	nonparam.	intervention target	al paired perfect inter- vention per node	model selection	$-\sum_k \log p_{m{ heta}}^k(\mathbf{x}^k))$	element-id. + graph-id	(c)
gel- et al. 21)	nonparam.	nonparam.	sample level on all real- izations of z_A^k	one imperfect fat- hand intervention	$\left\ g(\mathbf{x}^1)_{\hat{A}} - g(\mathbf{x}^2)_{\hat{A}}\right\ _2$	$-\sum_k H(g(\mathbf{x}^k)_{\hat{A}}), k \in \{1,2\}$	block-id.	

Under review as a conference paper at ICLR 2025

				2179 2180 2181 2182 2183	2170 2171 2172 2173 2174 2175 2176 2177 2178	2161 2162 2163 2164 2165 2166 2167 2168 2169	2154 2155 2156 2157 2158 2159 2160 2161	2152 2153
Work	Causal Model	Mixing Function	Invariance	Source of invari- ance, Inv. subset A	Invariance reg.	Sufficiency reg.	Identifiability Expl.	
Daunhawer et al. (2023)	nonparam.	nonparam.	sample level on all real- izations of z_A^k	one imperfect fat- hand intervention,	$\left\ g_1(\mathbf{x}^1)_{\hat{A}} - g_2(\mathbf{x}^2)_{\hat{A}}\right\ _2$	$-\sum_{k} H(g_{k}(\mathbf{x}^{k})_{\widehat{A}}),$ $k \in \{1,2\}$	block-id.	
Ahuja et al. (2022b)	nonparam.	nonparam.	sample level on all real- izations of z_A^k	one imperfect fat- hand intervention	$\left\ g(\mathbf{x}^1)_{\hat{A}} - g(\mathbf{x}^2)_{\hat{A}} + \delta\right\ _2$	$-\sum_k \mathbb{E}_{\mathbf{x}^k} \log p(\mathbf{x}^k g(\mathbf{x}^k)), \\ k \in \{1, 2\}$	block-id.	
Locatello et al. (2020)	nonparam.	nonparam.	sample level	one imperfect fat- hand intervention	avg. encoding	$-\sum_{k} \mathbb{E}_{\mathbf{x}^{k}} \log p(\mathbf{x}^{k} g(\mathbf{x}^{k})) ,$ $k \in \{1, 2\}$	block-id.	
Yao et al. (2023, Thm. 3.2)	nonparam.	nonparam.	sample level on all real- izations of z_A^k	partial observabil- ity	$\sum_{k,k' \in [K]} \sum_{\tilde{A}' = g_{k'}(\tilde{\mathbf{x}})_{\tilde{A}} \ _2$	$-\sum_{k\in[K]}H(g_k(\mathbf{x})_{\hat{A}})$	block-id.	
Yao et al. (2023, Thm. 3.8)	nonparam.	nonparam.	sample level on all real- izations of $z_{A_i}^k$	partial observabil- ity, $k \in V_i$	$\begin{split} & \sum_{k,k' \in V_i} \\ & \left\ g_k(\mathbf{x})_{\hat{A}(i,k)} - g_{k'}(\hat{\mathbf{x}})_{\hat{A}(i,k')} \right\ _2 \end{split}$	$-\sum_{k\in [K]} H(t_k\circ g_k(\mathbf{x}))$	block-id	
Brehmer et al. (2022)	nonparam.	nonparam.	sample level	perfect interven- tion per node	$\begin{array}{l} D_{\mathrm{KL}}\left(q(\mathcal{I},\hat{\mathbf{z}}^{1,2} \mid \mathbf{x}^{1,2}) \ p(\mathcal{I},\hat{\mathbf{z}}^{1,2}) \right) \\ \text{where } \hat{\mathbf{z}}^{k} := g(\mathbf{x}^{k}), k \in \{1,2\} \end{array}$	$) - \sum_k \mathbb{E}_{\mathbf{x}^k} \log p(\mathbf{x}^k g(\mathbf{x}^k)),$ $k \in \{1, 2\}$	element-id.	1
Lippe et al. (2022b)	nonparam.	nonparam.	transitional invari- ance on a distri- butional level	known-target inter- ventions x_t , invari- ant partition $A := [N] \setminus x_t$	$-H(\hat{s}_{t}^{t} \hat{z}^{t-1})$ where $\hat{z}^{t}:=g(\mathbf{x}^{t})$	$-p(\mathbf{x}^t \mathbf{x}^{t-1}, \mathcal{I}_t)$	block-id.	

~	ł	0	9
2	1	9	0
2	1	9	1
2	1	9	2
2	1	9	3
2	1	9	4
2	1	9	5
2	1	9	6
2	1	9	7
2	1	9	8
2	1	9	9
2	2	0	0
2	2	0	1
2	2	0	2
2	2	0	3
2	2	0	4
2	2	0	5
2	2	0	6
2	2	0	7
2	2	0	8
2	2	0	9
2	2	1	0
2	2	1	1
2	2	1	2

Work	Causal Model	Mixing Function	Invariance	Source of invari- ance, Inv. subset A	Invariance reg.	Sufficiency reg.	Identifiability	Expl.
Lippe et al. (2022a)	nonparam.	nonparam.	transitional invari- ance on a distri- butional level	known-target, par- tially perfect inter- ventions \mathcal{I}_t , invari- ant partition $A := [N] \setminus \mathcal{I}_t$	$-H(\hat{\mathbf{z}}_{A^t}^t \mid \hat{\mathbf{z}}^{t-1})$ where $\hat{\mathbf{z}}^t := g(\mathbf{x}^t)$	$-p(\mathbf{x}^t \mathbf{x}^{t-1},\mathcal{I}_t)$	block-id.	
Lippe et al. (2023)	nonparam.	nonparam.	transitional invari- ance on a distri- butional level	binary interven- tions (interven- tional target unknown)	$D_{\mathrm{KL}}(q(\mathbf{z}^t \mid \mathbf{x}^t) \parallel p(\mathbf{\hat{z}}^t \mid \mathbf{z}^{t-1}, \mathbf{r}^t))$ $\mathbf{r}^t \text{ observed regime vari-able}$, $-\log p(\mathbf{x}^{t} \mathbf{\hat{z}}^{t})$	block-id.	і 1
Lachapelle et al. (2023)	nonparam.	nonparam.	task sup- port	task distribution, overlapping task supports, number of causal variables known	$\sum_t \left\ \hat{\mathbf{w}}^{(t)} ight\ _{2,1}$	$\sum_t \mathcal{R}(\hat{\mathbf{w}}^{(t)} \circ_g)$	affine-id.	(e)
Fumero et al. (2024)	nonparam.	nonparam.	task sup- port	task distribution, overlapping task supports	$H(ilde{\mathbf{w}}){+}\sum_t \left\ \hat{\mathbf{w}}^{(t)} ight\ _1$	$\sum_t \mathcal{R}(\hat{\mathbf{w}}^{(t)} \circ_g)$	element-id.	(e)
Sagawa et al. (2019)	nonparam.	nonparam.	risk	invariant rela- tionship between label and invariant features, preserved under covariate shift	$\max_{k\in [K]} \mathcal{R}^k(w \circ g)$	$\max_{k \in [K]} \mathcal{R}^k(\mathbf{w} \circ g)$	NA	(f)

			2245 2246 2247 2248 2249	2236 2237 2238 2239 2240 2241 2242 2243 2243 2244	2227 2228 2229 2230 2231 2232 2233 2233 2234 2235	2222 2223 2224 2225 2226 2227	2218 2219 2220 2221	2217
Causal Model	Mixing Function	Invariance	Source of invari- ance, Inv. subset A	Invariance reg.	Sufficiency reg.	Identifiability	Expl.	
nonparam.	nonparam.	risk	invariant rela- tionship between label and invariant features, preserved under covariate shift	$\left\ \nabla_{\mathbf{w},\mathbf{w}=1} \mathcal{R}^k (\mathbf{w} \circ g) \right\ ^2$	$\sum_{k\in [K]} \mathcal{R}^k(\mathbf{w} \circ g)$	NA		
nonparam.	nonparam.	risk	invariant rela- tionship between label and invariant features, preserved under covariate shift	$\operatorname{Var}(\{\mathcal{R}^k(\mathbf{w} \circ g)\}_{k \in [K]})$	$\sum_{k\in [K]} \mathcal{R}^k(\mathbf{w} \circ g)$	NA	(f)	

Krueger et al. (2021)

Arjovsky et al. (2020)

Work

NA

 $\sum_{k \in [K]} \mathcal{R}^k(\mathbf{w} \circ g) \! + \! \operatorname{Var}(\mathcal{R})$

 $\left\| \nabla_{\mathbf{w},\mathbf{w}=1}\mathcal{R}^{k}\left(\mathbf{w}\circ g\right) \right\| ^{2}$

invariant rela-tionship between label and invariant features, preserved under covariate

risk

nonparam.

nonparam.

al. Ahuja et al. (2022a)

under shift