
Causal Effect Estimation with Mixed Latent Confounders and Post-treatment Variables

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Causal inference from observational data has attracted considerable attention among
2 researchers. One main obstacle is the handling of confounders. As direct mea-
3 surement of confounders may not be feasible, recent methods seek to address the
4 confounding bias via proxy variables, i.e., covariates postulated to be conducive to
5 the inference of latent confounders. However, the selected proxies may scramble
6 both confounders and post-treatment variables in practice, which risks biasing the
7 estimation by controlling for variables affected by the treatment. In this paper, we
8 systematically investigate the bias due to latent post-treatment variables, i.e., *latent*
9 *post-treatment bias*, in causal effect estimation. Specifically, we first derive the
10 bias when selected proxies scramble both confounders and post-treatment variables,
11 which we demonstrate can be arbitrarily bad. We then propose a novel Confounder-
12 identifiable VAE (CiVAE) to address the bias. Based on a mild assumption that the
13 prior of latent variables that generate the proxy belongs to a general exponential
14 family with at least one invertible sufficient statistic in the factorized part, CiVAE
15 *individually* identifies latent confounders and latent post-treatment variables up
16 to bijective transformations. We then prove that with individual identification,
17 the intractable disentanglement problem of latent confounders and post-treatment
18 variables can be transformed into a tractable independence test problem. Finally,
19 we prove that the true causal effects can be unbiasedly estimated with transformed
20 confounders inferred by CiVAE. Experiments on both simulated and real-world
21 datasets demonstrate significantly improved robustness of CiVAE.

22 1 Introduction

23 Causal inference, which aims to infer cause-and-effect relations from data, has gained increasing
24 prominence in various fields, such as social science, economics, and public health [10, 17, 34].
25 Traditional methods rely on the golden standard of randomized control trials (RCT) to draw valid
26 causal conclusions via experimentation [6]. Recently, more attention has been dedicated to causal
27 inference from observational data, where treatments, outcomes, and unit features are passively
28 observed, and researchers have no control over the treatment assignment mechanism [36, 37, 40].

29 One main obstacle to inferring valid causal relations from observational data is the confounding
30 bias, which occurs when we fail to account for the systematic difference between the treatment and
31 non-treatment group due to variables that causally influence the past treatments and the outcome, i.e.,
32 unobserved confounders [16]. If the confounders can be measured, a simple strategy to address the
33 bias is to control them via covariate adjustment [33] or propensity score re-weighting [24]. However,
34 confounders are not always measurable [23]. Therefore, recent methods seek to adjust for the
35 influence of unobserved confounders based on their proxies, which are easily acquirable covariates
36 postulated to be causally related with the unobserved confounders [29, 42, 28]. One exemplar work

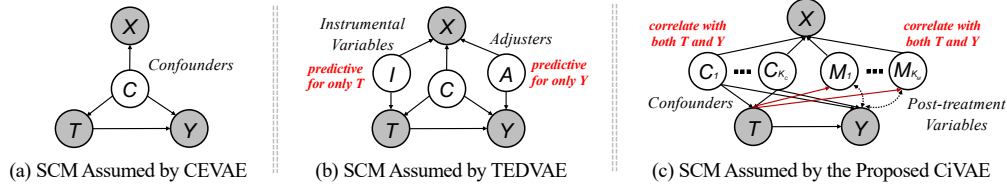


Figure 1: Comparison between the causal models assumed by CEVAE, TEDVAE, and CiVAE.

37 is the causal effect variational auto-encoder (CEVAE) [25], which has demonstrated that confounding
 38 bias can be mitigated by controlling latent variables inferred from the proxies of confounders.

39 Although proxy-based methods have achieved substantial progress in recent years, they may risk
 40 controlling latent post-treatment variables scrambled in the proxies, where **latent post-treatment**
 41 **bias** can be introduced. Here, we note that the negative effects of controlling *observed* post-treatment
 42 variables have been investigated in prior research [1, 9, 21]. For example, Montgomery et al. [30]
 43 found that more than 50% of the papers published in top journals of politics *inadvertently control*
 44 *post-treatment variables* in the experimental setting, even though researchers have complete control
 45 over which covariates to control for. On this basis, we postulate that the post-treatment bias could
 46 be even worse for proxy-based methods in the setting of observational study where variables are
 47 passively recorded. In addition, the post-treatment variables can be **latent** and scrambled into the
 48 observed covariates together with the latent confounders, which makes them difficult to disentangle.

49 Consider a real-world example from the Company¹. We found that *changing* a job from onsite to
 50 online mode causes applicants to make different decisions, and we want to estimate the causal effects
 51 of *switching a job from onsite to online mode to the decisions of the applicants* (reflected by statistics
 52 of applicants that apply for the job). In this case, the Company collected two groups of online (treated)
 53 and onsite (control) jobs, where the statistics of the applicants (e.g., the average age) are calculated as
 54 the surrogate outcome. Clearly, job seniority is a confounder, since less senior jobs are more likely to
 55 permit online work, and applicants for these jobs tend to be younger. However, the seniority level of
 56 a job can be difficult to measure. Therefore, the required skills of the job can be used as the proxy of
 57 the confounder "seniority", as senior jobs tend to require more advanced skills. However, a **caveat** is
 58 that switching to an online work mode may also alter the required skills of a job, thereby affecting the
 59 qualification and, therefore, the decision of the applicants. Consequently, directly using the skills as
 60 the proxy of the confounder "seniority" for adjustment could unintentionally control latent mediators
 61 (changed skills), which introduces latent post-treatment bias in the causal effect estimation.

62 Addressing the **latent post-treatment bias** faces multi-faceted challenges. First, there lacks a
 63 theoretical formulation of the bias when selected proxies scramble latent post-treatment variables
 64 for existing proxy-based methods. In addition, it is difficult to distinguish confounders and post-
 65 treatment variables in the latent space due to their similar observed behaviors. Existing covariate
 66 disentanglement-based methods, e.g., TEDVAE [44], focus on an easier task of disentangling latent
 67 confounders with latent adjusters and instrumental variables, which can be achieved by leveraging
 68 their different predictive abilities w.r.t. the treatment and outcome. However, since both latent
 69 confounders and post-treatment variables correlate with the treatment and the outcome, they cannot
 70 be disentangled by these methods. Finally, even if latent confounders can be distinguished from post-
 71 treatment variables, since most existing latent variable models have no identifiability guarantee [19],
 72 it is unclear whether controlling the inferred latent variables, which may be arbitrary transformations
 73 of the true confounders, can provide unbiased estimations of true causal effects.

74 To address the aforementioned challenges, we first analyze existing proxy-based methods when se-
 75 lected proxies scramble both latent confounders and post-treatment variables and show the estimation
 76 can be arbitrarily biased. We then propose a novel Confounder-identifiable VAE (CiVAE) to address
 77 the latent post-treatment bias. Specifically, we prove that based on a mild assumption that the prior
 78 of latent variables that generate the observed proxy (i.e., the latent confounders and post-treatment
 79 variables) belong to a general exponential family with at least one invertible sufficient statistic in the
 80 factorized part, latent confounders and latent post-treatment variables can be *individually* identified up
 81 to *simple bijective transformations*. With such identifiability guarantee, based on the causal relations
 82 among confounders, mediators, and treatment, we further demonstrate that the inferred confounders

¹Anonymized due to double-blind review policy.

83 (which are actually transformed proxies of the true confounders) could be properly distinguished
 84 from the latent post-treatment variables with pair-wise conditional independence tests. Finally, we
 85 prove that the true causal effects can be unbiasedly estimated based on transformed confounders
 86 inferred by CiVAE. Experiments on both simulated and real-world datasets demonstrate that CiVAE
 87 shows more robustness to latent post-treatment bias than existing methods.

88 2 Problem Formulation

89 In this paper, we assume the causal model in Fig. 1-(c). We use a binary random variable T to
 90 denote the treatment, a random vector $\mathbf{X} \in \mathbb{R}^{K_X}$ to denote the observed covariates (i.e., the proxy),
 91 and a random scalar $Y \in \mathbb{R}$ to denote the outcome. Furthermore, the observed covariates X are
 92 assumed to be generated from K_C independent latent confounders $\mathbf{C} \triangleq [C_1, C_2, \dots, C_{K_C}]$ causally
 93 influencing both T and Y , and K_M latent post-treatment variables $\mathbf{M} \triangleq [M_1, M_2, \dots, M_{K_M}]$ under
 94 the causal influence of the treatment (where the relation between \mathbf{M} and Y can be arbitrary). We use
 95 the random vector $\mathbf{Z} \triangleq [\mathbf{C} || \mathbf{M}] \in \mathbb{R}^{K_Z=K_C+K_M}$ to denote all latent factors. **Our aim** is to estimate
 96 the average causal effects of treatment T on outcome Y with auxiliary confounder information in \mathbf{X} ,
 97 where the estimation should be devoid of both confounding bias and post-treatment bias.

98 3 Theoretical Analysis of Latent Post-Treatment Bias

99 3.1 Preliminaries and Assumptions

100 To achieve such a purpose, we first define the (conditional) average treatment effects (C/ATE) when
 101 covariates \mathbf{X} scramble both latent confounders \mathbf{C} and post-treatment variables \mathbf{M} . We then define
 102 the post-treatment bias when covariates \mathbf{X} are directly used as the proxy of confounders. To facilitate
 103 the analysis, we make the following assumption regarding the causal generative process.

104 **Assumption 1. (Noisy-Injectivity).** We assume $\mathbf{X} = f(\mathbf{C}, \mathbf{M}) + \epsilon$, where f is a deterministic
 105 function that combines latent confounders \mathbf{C} and latent post-treatment variables \mathbf{M} into observations
 106 \mathbf{X} , and ϵ is random noise. In addition, we assume that the function f is **injective**; beyond injectivity,
 107 f can be arbitrarily nonlinear. We use $f^\dagger : \mathbf{X} \rightarrow [\mathbf{C} || \mathbf{M}]$ to denote its left inverse. We use
 108 $f_C^\dagger : \mathbf{X} \rightarrow \mathbf{C}$ and $f_M^\dagger : \mathbf{X} \rightarrow \mathbf{M}$ to denote the mapping from \mathbf{X} to \mathbf{C} , \mathbf{M} , respectively.

109 **Noisy-Injectivity** is a common assumption made either explicitly or implicitly in most existing proxy-
 110 of-confounder-based causal inference algorithms. For example, if both \mathbf{X} and \mathbf{C} are categorical,
 111 [31] assumes that \mathbf{X} has at least the same number of categories as \mathbf{C} , whereas the effect restoration
 112 algorithm [35] assumes that the matrix of $p(\mathbf{C}, \mathbf{X})$ to be full-rank. Although CEVAE [25] makes no
 113 explicit injectivity assumption between \mathbf{C} and \mathbf{X} , it requires that the joint distribution $p(\mathbf{C}, \mathbf{X}, T, Y)$
 114 can be fully recovered from the observations (\mathbf{X}, T, Y) . [2] show that some of the possible identifica-
 115 tion criteria for the recovery include **1)** having multiple independent views of \mathbf{C} in \mathbf{X} [8], and **2)** \mathbf{C}
 116 is categorical and \mathbf{X} is a mixture of Gaussian components determined by \mathbf{C} (that is, \mathbf{X} is generated
 117 by bijective mapping of \mathbf{C} to the mean of the corresponding component with added Gaussian noise).

118 In the following part of this section, we omit the noise ϵ to gain better intuition of latent post-treatment
 119 bias (but all the exact conclusions will still hold in the posterior sense [19]). In Section 4, we assume
 120 noise exists and demonstrate that our method can still properly identify the latent confounders.

121 3.2 Causal Estimand and the True ATE

122 Based on Assumption 1, we are ready to define the estimand of average treatment effect (ATE)
 123 through controlling the covariates \mathbf{X}' , as well the as the true (conditional) average treatment effects.

124 **Definition 1. (DCEV & DEV).** We define the Difference in Conditional Expected Values (DCEV) as:

$$DCEV(\mathbf{x}') = \mathbb{E}[Y|T = 1, \mathbf{X}' = \mathbf{x}'] - \mathbb{E}[Y|T = 0, \mathbf{X}' = \mathbf{x}'], \quad (1)$$

125 which is the difference of the expected value of Y for units with variable $\mathbf{X}' = \mathbf{x}'$ in the treatment
 126 group and the non-treatment group. Based on $DCEV(\mathbf{x}')$, we define the Difference in Expected
 127 Value (DEV) as $DEV(\mathbf{X}') = \mathbb{E}_{p(\mathbf{X}')} [DCEV(\mathbf{X}')] as the expectation of DCEV w.r.t. $p(\mathbf{X}')$.$

128 $DEV(\mathbf{X}')$ denotes the estimand of ATE when \mathbf{X}' is the covariates that we choose to control (i.e.,
 129 calculate the expected difference in each stratum of $\mathbf{X}' = \mathbf{x}'$). If $\mathbf{X}' = \emptyset$, $DEV(\emptyset)$ represents
 130 the *naive estimator* that directly calculates the expected difference of the outcome Y between the
 131 treatment group and the non-treatment group. With the causal estimand $DEV(\mathbf{X}')$ defined, we then
 132 derive the true causal effects with the covariates \mathbf{X}' when it scrambles both latent confounders and
 133 post-treatment variables according to the generative process described in Assumption 1:

134 **Definition 2.** Under Assumption 1, we define the *Conditional Average Treatment Effect (CATE)* for
 135 individuals with observed covariates $\mathbf{X} = \mathbf{x}$ by controlling only the confounder part in \mathbf{X} as:

$$CATE(\mathbf{x}) = \mathbb{E}[Y|T = 1, \mathbf{C} = f_C^\dagger(\mathbf{x})] - \mathbb{E}[Y|T = 0, \mathbf{C} = f_C^\dagger(\mathbf{x})], \quad (2)$$

136 with the *Average Treatment Effect (ATE)* of treatment T defined as:

$$ATE = \mathbb{E}[Y|do(T = 1)] - \mathbb{E}[Y|do(T = 0)] = \mathbb{E}_{p(\mathbf{C})}[\mathbb{E}[Y|T = 1, \mathbf{C}] - \mathbb{E}[Y|T = 0, \mathbf{C}]]. \quad (3)$$

137 Please note that we only consider the latent confounder component of the observed features \mathbf{X} in the
 138 definition of CATE in Eq. (2). This is because the causal relationship between the post-treatment
 139 variables \mathbf{M} and the outcome Y is indeterminate. However, if the specific relationship between \mathbf{M}
 140 and Y can be further established by the researcher (e.g., all elements of \mathbf{M} are latent mediators),
 141 more precise forms of CATE can be derived with path-specific counterfactual analysis [5, 14].

142 3.3 Latent Post-Treatment Bias

143 With $DEV(\mathbf{X}')$ (the ATE estimator that control for the covariates \mathbf{X}'), CATE, and ATE defined in
 144 Section 3.2, in this section, we analyze the *latent post-treatment bias* of existing proxy-of-confounder-
 145 based causal inference methods, such as CEVAE, that control for latent variables inferred from
 146 the covariates \mathbf{X} to estimate the ATE of T on Y , when \mathbf{X} scrambles both latent confounders and
 147 post-treatment variables as Assumption 1. In our analysis, Lemma 3.1 will be frequently used.

148 **Lemma 3.1.** For an injective function g , $\mathbb{E}[Y|\mathbf{X}' = \mathbf{x}'] = \mathbb{E}[Y|g(\mathbf{X}') = g(\mathbf{x}')] holds.$

149 The proof when g is differentiable *a.e.* can be referred to in Appendix C.1. Since the latent variable
 150 models used in existing methods (such as VAE with factorized Gaussian prior in CEVAE) lack
 151 identifiability guarantee (i.e., the recovery of the exact latent variables), we assume that these models
 152 can recover the true latent space $\mathbf{Z} = [\mathbf{C}, \mathbf{M}]$ up to invertible transformations \tilde{f} , where the inference
 153 process can be represented as $\tilde{\mathbf{Z}} = \tilde{f}(\mathbf{X}) = \tilde{f} \circ f^\dagger(\mathbf{X})$. With such an assumption, we have the
 154 following theorem regarding the latent post-treatment bias when \mathbf{X} mixes post-treatment variables.

155 **Theorem 3.2.** If the observed covariates \mathbf{X} are generated from latent confounders \mathbf{C} and latent
 156 post-treatment variables \mathbf{M} according to Assumption 1, the latent post-treatment bias of a proxy-
 157 based causal inference algorithm that controls latent variables $\tilde{\mathbf{Z}}$ inferred from \mathbf{X} via $\tilde{f} = \tilde{f} \circ f^\dagger$:
 158 $\mathbb{R}^{K_X} \rightarrow \mathbb{R}^{K_C + K_M}$ to estimate the ATE can be formulated as follows:

$$\begin{aligned} Bias(\mathbf{X}) &= ATE - DEV(\tilde{f}(\mathbf{X})) = ATE - \mathbb{E}[\mathbb{E}[Y|T = 1, \tilde{f}(\mathbf{X})] - \mathbb{E}[Y|T = 0, \tilde{f}(\mathbf{X})]] \\ &= ATE - \mathbb{E}[\mathbb{E}[Y|1, \tilde{f} \circ f^\dagger(f(\mathbf{C}, \mathbf{M}))] - \mathbb{E}[Y|0, \tilde{f} \circ f^\dagger(f(\mathbf{C}, \mathbf{M}))]] \\ &= \mathbb{E}[\mathbb{E}[Y|1, \mathbf{C}] - \mathbb{E}[Y|0, \mathbf{C}]] - \mathbb{E}[\mathbb{E}[Y|1, \mathbf{C}, \mathbf{M}] - \mathbb{E}[Y|0, \mathbf{C}, \mathbf{M}]], \end{aligned} \quad (4)$$

159 which can be arbitrarily bad. Therefore, the estimator of existing proxy-of-confounder-based meth-
 160 ods, i.e., $DEV(\tilde{f}(\mathbf{X}))$, is an arbitrarily biased estimator of the ATE, when the selected proxy of
 161 confounders \mathbf{X} accidentally mixes in latent post-treatment variables \mathbf{M} .

162 The final step of Eq. (4) can be proved since f is injective and \tilde{f} bijective, the composite $\tilde{f} \circ f^\dagger \circ f$:
 163 $[\mathbf{C}, \mathbf{M}] \rightarrow \tilde{\mathbf{Z}}$ is bijective, so we can use Lemma 3.1 to remove $\tilde{f} \circ f^\dagger \circ f$ in the condition.

164 3.4 Examples in the Linear Case

165 Generally, the latent post-treatment bias defined in Eq. (4) cannot be simplified, because (i) the
 166 causal relationship between \mathbf{M} and Y are indeterminate, and (ii) the causal influence of \mathbf{C} , \mathbf{M} ,
 167 and T on Y can be arbitrary. However, for linear structural causal models with determined causal
 168 relationships between \mathbf{M} and Y (e.g., \mathbf{M} are mediators, which are post-treatment variables that have
 169 causal influences on the outcomes), stronger conclusions can be drawn as follows:

170 **Corollary 3.3. (Mixed Latent Mediator).** For the linear Structural Causal Model (SCM) defined as:

$$\begin{aligned}
& (i) T \leftarrow \mathbb{1}(\alpha_T + \sum \beta_i \cdot C_i > a), \quad (ii) M_j \leftarrow \alpha_M + \gamma_j \cdot T \\
& (iii) \mathbf{X} \leftarrow \alpha_X + \mathbf{A}[\mathbf{M}|\mathbf{C}], \quad (iv) Y \leftarrow \alpha_Y + \tau \cdot T + \sum \theta_j \cdot M_j + \sum \kappa_i \cdot C_i,
\end{aligned} \tag{5}$$

171 where the mixture function $f = \mathbf{A} \in \mathbb{R}^{K_X \times (K_C + K_M)}$ is a full column-rank matrix, the CATE, ATE,
172 and the bias of proxy-of-confounder-based causal inference model that controls the latent variables
173 $\hat{\mathbf{Z}}$ inferred via $\hat{\mathbf{Z}} = \tilde{f}(\mathbf{X}) = \mathbf{B}^T \mathbf{X}$ can be formulated as follows:

$$\begin{aligned}
& ATE = CATE = \tau + \sum \gamma_j \cdot \theta_j, \quad \text{and } DEV(\hat{\mathbf{Z}}) = \mathbb{E}[DCEV(\hat{\mathbf{Z}})] = DCEV(\hat{\mathbf{Z}}) = \tau \\
& Bias(\hat{\mathbf{Z}}) = ATE - DEV(\hat{\mathbf{Z}}) = \sum \gamma_j \cdot \theta_j,
\end{aligned} \tag{6}$$

174 where $\mathbf{B} \in \mathbb{R}^{K_X \times (K_C + K_M)}$ is another full column-rank matrix. Since $\sum \gamma_j \cdot \theta_j$ is arbitrary, the
175 estimator $DEV(\hat{\mathbf{Z}}) = \mathbb{E}[DCEV(\mathbf{B}^T \mathbf{X})]$ is arbitrarily biased for ATE estimation.

176 The proof of Eq. (6) is provided in Appendix C.2. In addition, we show that post-treatment variables
177 \mathbf{M} DO NOT necessarily need to have direct causal effects on the outcome Y to incur arbitrary bias
178 in ATE estimation. In Appendix C.3, we provide another example (i.e., Mixed Latent Correlator) in
179 the linear case where \mathbf{M} is correlated with Y through unobserved confounders U in Corollary C.1.

180 4 Methodology

181 In this section, we introduce the proposed Confounder-identifiable Variational Auto-Encoder (CiVAE)
182 in detail. Specifically, we first prove that if the prior distribution of the true latent variables $\mathbf{Z} =$
183 $[\mathbf{C}, \mathbf{M}]$ satisfies certain weak assumptions, CiVAE *individually* identify $[\mathbf{C}, \mathbf{M}]$ up to bijective
184 transformations. Then, utilizing the causal relations between \mathbf{C} , \mathbf{M} , and T , we novelly transform the
185 challenging confounder-identifiability problem into a tractable pair-wise conditional independence
186 test problem, which can be effectively solved with kernel-based methods. The generalization of
187 CiVAE to address the interactions among $[\mathbf{C}, \mathbf{M}]$ are discussed in Section D of the Appendix.

188 4.1 Generative Process

189 The fundamental work on the identifiability of deep variational inference, i.e., the identifiable VAE
190 (iVAE) [19], makes a strict assumption that the prior of true latent variables \mathbf{Z} (i.e., $[\mathbf{C}, \mathbf{M}]$ in
191 our case) is conditionally factorized given the available covariates. However, since both \mathbf{C} and
192 \mathbf{M} form fork structures with the outcome Y (see Fig. 1-(c)) [22], $C_i, C_j, M_i,$ and M_j are not
193 independent given Y . Recently, Non-Factorized iVAE (NF-iVAE) [26] was proposed that allows
194 arbitrary dependence among the true latent variables \mathbf{Z} in the conditional priors, where \mathbf{Z} can be
195 identified up to arbitrary non-linear transformations. However, the transformation is not necessarily
196 invertible, which is risky as multiple values of the confounders may collapse, leading to bias when
197 estimating the ATE by averaging the $DCEV$ calculated in each stratum of the inferred confounders.

198 In contrast to NF-iVAE, CiVAE guarantees the individual and bijective identifiability of \mathbf{Z} by putting
199 a general exponential family *with at least one invertible sufficient statistic in the factorized part* as its
200 prior when conditioning on treatment T and outcome Y , which can be formulated as follows.

201 **Assumption 2.** Let $\mathbf{Z} = [\mathbf{C}|\mathbf{M}]$ be the random vector for latent variables that causally gener-
202 ate the observed covariates \mathbf{X} according to Assumption 1. We assume that the conditional
203 prior of \mathbf{Z} given the outcome Y and the treatment T belongs to a general exponential family
204 with parameter vector $\boldsymbol{\lambda}(Y, T)$ and sufficient statistics $\mathbf{S}(\mathbf{Z}) = [\mathbf{S}_f(\mathbf{Z})^T, \mathbf{S}_{nf}(\mathbf{Z})^T]^T$. Specif-
205 ically, $\mathbf{S}(\mathbf{Z})$ is composed of (i) the sufficient statistics of a factorized exponential family, i.e.,
206 $\mathbf{S}_f(\mathbf{Z}) = [\mathbf{S}_1(Z_1)^T, \dots, \mathbf{S}_{K_Z}(Z_{K_Z})^T]^T$, where all components $\mathbf{S}_i(Z_i)$ have dimension larger
207 than or equal to 2 and *each \mathbf{S}_i has at least one invertible dimension*, and (ii) $\mathbf{S}_{nf}(\mathbf{Z})$, where \mathbf{S}_{nf} is
208 a neural network with ReLU activation. The density of the conditional prior can be formulated as:

$$p_{\mathbf{S}, \boldsymbol{\lambda}}(\mathbf{Z}|Y, T) = \mathcal{Q}(\mathbf{Z}) / \mathcal{C}(Y, T) \exp[\mathbf{S}(\mathbf{Z})^T \boldsymbol{\lambda}(Y, T)], \tag{7}$$

209 where $\mathcal{Q}(\mathbf{Z})$ is the base measure, and $\mathcal{C}(Y, T)$ is the normalizing constant independent of \mathbf{Z} .

210 We justify that assumption 2 is weak and practical as follows. (i) Neural networks with ReLU
 211 activation have **universal approximation ability** of distributions [27]. Therefore, Eq. (7) can model
 212 arbitrary dependence between true latent confounders \mathbf{C} and post-treatment variables \mathbf{M} conditional
 213 on T and Y . (ii) Although CiVAE makes an extra assumption that $\forall i$, at least one dimension of \mathbf{S}_i is
 214 invertible, this can be easily satisfied as most commonly used exponential family distributions, such
 215 as Gaussian, Bernoulli, etc., has at least one invertible sufficient statistics².

216 The reason why we use ReLU as the activation is that, the identifiability of iVAE relies on the
 217 condition that the sufficient statistics \mathbf{S} have zero second-order cross-derivative. The factorized part,
 218 i.e., \mathbf{S}_f , satisfies it trivially as all cross-derivatives of \mathbf{S}_f are zero. In addition, since the ReLU neural
 219 networks are linear *a.e.*, all second-order derivatives of $\mathbf{S}_{n,f}$ are zero. Therefore, identifiability holds
 220 after adding $\mathbf{S}_{n,f}$ in the prior that allows the capturing of arbitrary dependence among \mathbf{Z} .

221 4.2 Optimization Objective

222 Combining Assumptions 1 and 2, the generative process assumed by CiVAE can be formulated as:

$$(i) p_{\theta}(\mathbf{X}, \mathbf{Z} | Y, T) = p_f(\mathbf{X} | \mathbf{Z}), (ii) p_{\mathbf{S}, \lambda}(\mathbf{Z} | Y, T), (iii) p_f(\mathbf{X} | \mathbf{Z}) = p_{\epsilon}(\mathbf{X} - f(\mathbf{Z})). \quad (8)$$

223 where $\theta = (f, \lambda, \mathbf{S}) \in \Theta$ are the parameters of the generative distribution. Since the generative
 224 process of CiVAE is parameterized by deep neural networks, the posterior distribution of \mathbf{Z} , i.e.,
 225 $p_{\theta}(\mathbf{Z} | \mathbf{X}, Y, T)$, is intractable. Therefore, we resort to variational inference [4], where we introduce
 226 an approximate posterior $q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T)$ parameterized by a deep neural network with a trainable
 227 parameter ϕ , and in $q_{\phi}(\mathbf{Z} | \cdot)$ finds the one closest to $p_{\theta}(\mathbf{Z} | \cdot)$ measured by KL divergence. The
 228 minimization of KL is equivalent to maximization of the evidence lower bound (ELBO):

$$\mathcal{L}(\theta, \phi) := \mathbb{E}_{q_{\phi}} \left[\log p_f(\mathbf{X} | \mathbf{Z}) + \underbrace{\log p_{\mathbf{S}, \lambda}(\mathbf{Z} | Y, T) - \log q_{\phi}(\mathbf{Z} | \cdot)}_{\text{KL of posterior with prior}} \right]. \quad (9)$$

229 Since the normalization constant \mathcal{C} in Eq. (7) is generally intractable, it is infeasible to directly learn
 230 \mathbf{S}, λ by optimizing Eq. (9). Therefore, we substitute the KL term in Eq. (9) with the widely-used
 231 score matching [13] to learn unnormalized densities instead as follows:

$$\begin{aligned} \mathcal{L}(\mathbf{S}, \lambda, \phi) &:= \mathbb{E}_{q_{\phi}(\mathbf{Z} | \cdot)} \left[\|\nabla_{\mathbf{Z}} \log q_{\phi}(\mathbf{Z} | \cdot) - \nabla_{\mathbf{Z}} \log p_{\mathbf{S}, \lambda}(\mathbf{Z} | Y, T)\|^2 \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{Z} | \cdot)} \left[\sum_{j=1}^{K_{\mathbf{Z}}} \left[\frac{\partial^2 p_{\mathbf{S}, \lambda}(\mathbf{Z} | Y, T)}{\partial Z_j^2} + \frac{1}{2} \left(\frac{\partial p_{\mathbf{S}, \lambda}(\mathbf{Z} | Y, T)}{\partial Z_j} \right)^2 \right] \right] + \text{const}, \end{aligned} \quad (10)$$

232 4.3 Identifiability of CiVAE

233 With the generative process and optimization objective of CiVAE discussed in previous sub-sections,
 234 we are ready to introduce the final assumption of CiVAE, which, combined with Assumptions 1 and
 235 2, leads to the main Theorem of this paper, which states the identifiability of CiVAE.

236 **Assumption 3.** Assume the following: (i) The set $\{\mathbf{X} \in \mathcal{X} | \phi(\mathbf{X}) = 0\}$ has measure zero, where ϕ
 237 is the characteristic function of the density p_f in Eq. (8). (ii) The sufficient statistics, \mathbf{S}_i in \mathbf{S}_f are all
 238 twice differentiable. (iii) The mixture function f in Eq. (8) has all second-order cross derivatives.
 239 (iv) There exist $k + 1$ distinct points $(Y, T)_0, \dots, (Y, T)_k$ s.t. the matrix $\mathbf{L} = [\lambda((Y, T)_1) -$
 240 $\lambda((Y, T)_0), \dots, \lambda((Y, T)_k) - \lambda((Y, T)_0)]$ of size $k \times k$ is invertible, where $k = \text{Dim}(\mathbf{S})$.

241 Here, we note that Assumptions (i) - (iii) are trivial for differentiable neural networks. The Assumption
 242 (iv) can be intuitively understood as independent samples of (Y, T) are required to identify \mathbf{C} and
 243 \mathbf{M} . The identifiability theorem of CiVAE can be formulated as follows.

244 **Theorem 4.1.** If Assumptions 1, 2, and 3 hold, and if $\theta, \tilde{\theta} \in \Theta \rightarrow p_{\theta}(\mathbf{X} | Y, T) = p_{\tilde{\theta}}(\mathbf{X} | Y, T)$, the
 245 true latent variables \mathbf{Z} are identifiable up to **permutation** and **element-wise bijective transformation**.
 246 Furthermore, in the case of **variational inference**, if we denote the true parameter that generates the
 247 data as θ^* , if (i) the distribution family $q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T)$ contains the posterior $p_{\theta}(\mathbf{Z} | \mathbf{X}, Y, T)$, and
 248 $q_{\phi}(\mathbf{Z} | \mathbf{X}, Y, T) > 0$, (ii) we optimize Eq. (4) w.r.t. both θ, ϕ , then in the limit of infinite data, true
 249 parameters θ^* can be learned up to a permutation and bijective transformation of \mathbf{Z} .

²There are a few exponential family dist. with no invertible sufficient statistics, e.g., Weibull with even shape parameter k . However, these distributions are not commonly used in statistics or machine learning.

250 The proof of Theorem 4.1 non trivially extends the NF-iVAE paper [26] by incorporating the new
 251 assumption introduced in CiVAE (i.e., each S_i has at least one invertible dimension) to ensure that the
 252 transformation of each Z_i is bijective. The detailed proof is provided in Appendix C.4 for reference.

253 4.4 Identification of Latent Confounders

254 Theorem 4.1 ensures that the latent variables \hat{Z} inferred by CiVAE cannot **(i)** mix confounders
 255 and post-treatment variables in each dimension, or **(ii)** collapsing of different values of the latent
 256 confounders into the same value. To further determine the dimensions of confounder and post-
 257 treatment variable in \hat{Z} , we rely on the causal relations between latent variables \hat{Z} and the treatment
 258 T and the associated marginal/conditional dependence properties, which are discussed as follows.

- 259 • **Case 1. Intra-Confounders.** Latent confounders C_i, C_j and the treatment T form the V structure
 260 $C_i \rightarrow T \leftarrow C_j$. Therefore, C_i and C_j are marginally **independent**, whereas they become
 261 **dependent** when conditioning on the assigned treatment T .
- 262 • **Case 2. Intra-Post Treatment Variables.** Latent post-treatment variables M_i, M_j and the treatment
 263 T form a *Fork-structure* $M_i \leftarrow T \rightarrow M_j$, where M_i, M_j are marginally **dependent**, but they
 264 become **independent** after conditioning on the assigned treatment T .
- 265 • **Case 3. Cross-Confounder and Post-Treatment Variables.** Latent confounder C_i , latent post-
 266 treatment variable M_j , and the treatment T forms a Chain structure $C_i \rightarrow T \rightarrow M_j$, where $C_i,$
 267 M_j are marginally dependent, and they become **independent** after conditioning on T .

268 From the above analysis we can find that, the dependence between two latent variables \hat{Z}_i and \hat{Z}_j
 269 **increases** after conditioning on the treatment T ONLY in the case of *intra-confounders*. Therefore,
 270 if more than one latent confounder exists, which is highly probable when covariates \mathbf{X} are high-
 271 dimensional, we can conduct independence test $\text{Ind}(\hat{Z}_i, \hat{Z}_j)$ and $\text{CInd}(\hat{Z}_i, \hat{Z}_j|T)$ for all pairs of
 272 inferred latent variables, which can be implemented via kernel-based methods as [43], and select
 273 the pairs where the p-value of CInd is larger than that of Ind as latent confounders. Here, we note
 274 that the kernel-based (conditional) independence test incurs $N^2 \times K_Z^2$ complexity in the training
 275 phase. However, once the dimensions of the confounders in \hat{Z} are determined, CiVAE **has the same**
 276 **complexity as CEVAE** for the estimation of CATE and ATE in the test phase.

277 4.5 ATE Estimator with Transformed Confounders

278 Finally, we demonstrate that controlling the transformed confounders \hat{C} inferred by CiVAE provides
 279 an unbiased estimation of ATE. Specifically, we have the final Theorem show the unbiasedness.

280 **Theorem 4.2.** *Controlling bijective of confounders is equivalent to original confounders in ATE*
 281 *estimation, i.e., $DEV(\hat{C}) = DEV(g(C)) = ATE$, if the transformation function g is bijective.*

282 The proof of Theorem 4.2 for discrete C is trivial (where $\hat{C} = g(C)$ represents a simple relabeling
 283 of the stratum that we calculate the $DCEV$ and take the expectation). The proof in the continuous
 284 case where g is differentiable is provided in Appendix C.5. With Theorem 4.2, we can control the
 285 identified latent confounders as true confounders, providing an unbiased estimate of ATE.

286 5 Empirical Study

287 In this section, we provide and analyze the experiments we conduct on both simulated and real-world
 288 datasets, where a code demo written in PyTorch and Pyro is provided in this anonymous URL.

289 5.1 Datasets

290 **Simulated Datasets.** We first establish two simulated datasets, i.e., `LatentMediator` and
 291 `LatentCorrelator`, that consider two types of post-treatment variables, i.e., **(i)** mediators and
 292 **(ii)** correlators, i.e., variables that are correlated with the outcome Y via latent confounders U , where
 293 the causal generative process is under the full control of the experimenter. The generative process of
 294 the two datasets can be referred to in Corollary 3.3 and Corollary C.1 in the Appendix, respectively.
 295 In our experiments, C are generated from Gaussian distribution as $C \sim \text{Gaussian}(0, \mathbf{I}_{K_C})$. For

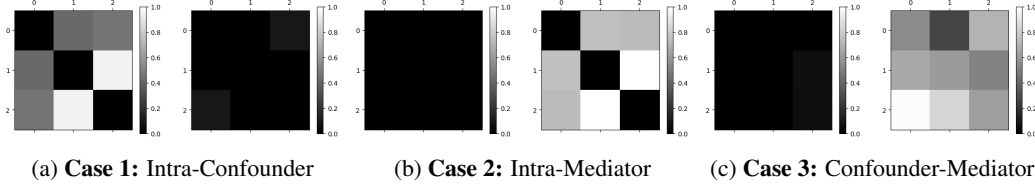


Figure 2: Visualization of p -value of independence test before and after conditioning on treatment T .

296 LatentMediator, γ is set as $[-1, -1, -1]$, θ is set as $[1, 1, 1]$, and τ is set as 2, which results in
 297 $ATE = -1$. For the LatentCorrelator dataset, we set the same γ and θ as the LatentMediator
 298 dataset, where parameters ϕ and τ are set to 1, which results in an overall ATE of 1.

299 **Real-world Datasets.** In addition, we build real-world datasets from the Company to estimate the
 300 ATE of *switching a job from onsite to online work mode to the statistics of the applicants*. The
 301 average age and the variance of gender of the applicants are two outcomes of interest. Covariates
 302 $\mathbf{X} \in \{0, 1\}^{K_x}$ include the required skills of the job. Specifically, we establish a cohort of 3,228
 303 jobs from the Bay Area in the US, where a preliminary study shows that $DEV(\emptyset) \approx 2$ years³ (i.e.,
 304 online job applicants are two years younger than onsite job applicants in the collected data), and
 305 $DEV(\emptyset) \approx -0.015$ (i.e., online jobs exhibit 0.015 more gender variance than onsite jobs in the
 306 collected data). To simulate \mathbf{C} and \mathbf{M} , we first learn a generative model as follows:

$$\mathbf{Z} \sim \text{Gaussian}(\mathbf{0}, \mathbf{I}_{K_Z}), \mathbf{X} \sim \text{Multi}(NN_f(\mathbf{Z})), Y \sim \text{Gaussian}(\mathbf{w} \odot \mathbf{Z}, 1), \quad (11)$$

307 where Multi represents multinomial distribution, NN_f is a neural network with softmax activation,
 308 $\mathbf{Z}, \mathbf{w} \in \mathbb{R}^{K_Z}$, $K_Z = 8$, and \odot represents the element-wise product operator, respectively. We
 309 then treat the first $K_C = 5$ dimensions of \mathbf{Z} as the latent confounders \mathbf{C} and the remaining
 310 $K_M = K_Z - K_C$ dimensions as the latent mediators \mathbf{M} . After learning NN_f and \mathbf{w} according to
 311 Eq. (11), we draw latent confounders $\mathbf{C} \in \text{Gaussian}(0, \mathbf{I})$, latent mediators $\mathbf{M} = T \cdot \gamma$, and set the
 312 outcome $Y = \mathbf{w} \odot [\mathbf{C} || \mathbf{M}] + \tau \cdot T$, where the true ATE can be calculated as $\text{sum}(\gamma \odot \mathbf{w}_{-K_M}) + \tau$.

313 5.1.1 Disentangle Confounders and Post-treatment Variables

314 We first show the p -value of the kernel-based pairwise independence test of the true latent variables
 315 before and after conditioning on the assigned treatment T . From Fig. 2, we can find that the distinction
 316 of the intra-confounder case from the other two cases discussed in Subsection 4.4 is significant. Here,
 317 we should note this relies on the assumption that latent confounders are independent. If the latent
 318 confounders are correlated, we can first use causal discovery techniques such as the PC algorithm [39]
 319 to find direct parents of T , and use our algorithm as the refinement to determine the true confounders
 320 \mathbf{C} from the misidentified post-treatment variables (Experiments see Section D) in Appendix.

321 5.2 Baselines

322 The baselines we include for comparisons can be categorized into three classes. **(i) Unawareness**,
 323 where no information in \mathbf{X} is used for ATE estimation. We implement the naive LR0 estimator, which
 324 regresses Y on T and uses the coefficient to estimate the ATE [15] (LR0 equals to $DEV(\emptyset)$, i.e., the
 325 difference of the average outcome between the treatment and non-treatment group). **(ii) Control- \mathbf{X}** ,
 326 which directly controls the covariates \mathbf{X} . In this class, LR1 regresses Y on T and \mathbf{X} , whereas TarNet
 327 uses a two-branch neural network to estimate the $DEV(\mathbf{X})$ **(iii) Control- \mathbf{Z}** , which controls latent
 328 variables \mathbf{Z} learned from the covariates \mathbf{X} . Methods from this class include the CEVAE [25] and
 329 covariate disentanglement methods, such as DR-CFR [12], TEDVAE [44], NICE [38], and AFS [41].

330 5.2.1 Results and Analysis

331 From Table 1, we can find that for all four datasets, CEVAE is worse than the naive LR0 estimator.
 332 In addition, for the LatentMediator and Company (Age) dataset, all methods except CiVAE fail
 333 to predict the negativity of the ATE . Covariates disentanglement-based methods, i.e., DR-CFR
 334 and TEDVAE, inherit the latent post-treatment bias of CEVAE. The reason is that, these methods
 335 disentangle latent confounders \mathbf{C} from latent instrumental variables \mathbf{I} and latent adjusters \mathbf{A} by

³which leads to 0.178 and -0.105 after standardization of the outcome.

Table 1: Comparison of CiVAE with baselines under latent post-treatment bias on various datasets.

Dataset	LatentMediator		LatentCorrelator		Company (Age)		Company (Gender)	
	ATE.	Err.	ATE.	Err.	ATE.	Err.	ATE.	Err.
LR0	0.975 ± 0.032	1.975	2.977 ± 0.032	1.977	0.131 ± 0.015	0.399	-0.105 ± 0.009	-0.213
LR1	1.457 ± 0.167	2.457	3.400 ± 0.130	2.400	0.093 ± 0.029	0.361	-0.175 ± 0.014	-0.256
TarNet	1.461 ± 0.172	2.461	3.414 ± 0.146	2.414	0.112 ± 0.085	0.380	-0.167 ± 0.021	-0.248
CEVAE	1.550 ± 0.292	2.550	3.323 ± 0.167	2.323	0.106 ± 0.078	0.374	-0.180 ± 0.028	-0.261
DR-CFR	1.239 ± 0.324	2.239	3.185 ± 0.319	2.185	0.094 ± 0.089	0.362	-0.159 ± 0.030	-0.240
NICE	1.868 ± 0.530	2.868	1.942 ± 0.524	0.942	0.149 ± 0.126	0.417	-0.186 ± 0.041	-0.267
TEDVAE	1.042 ± 0.315	2.042	3.138 ± 0.281	2.138	0.097 ± 0.093	0.365	-0.143 ± 0.027	-0.224
AFS	1.496 ± 0.825	2.496	3.251 ± 0.398	2.251	0.105 ± 0.102	0.373	-0.163 ± 0.045	-0.244
CiVAE	-0.822 ± 0.753	0.178	1.199 ± 0.765	0.199	-0.140 ± 0.137	0.128	-0.106 ± 0.064	-0.187
True ATE	-1.000 ± 0.000	0.000	1.000 ± 0.000	0.000	-0.268 ± 0.000	0.000	-0.081 ± 0.000	0.000

utilizing their causal relations with T and Y , i.e., I is predictive only for T , A is predictive only for Y , whereas C is predictive for both T and Y . For example, TEDVAE includes three encoders to infer three sets of latent variables $\hat{I}, \hat{A}, \hat{C}$ from X and adds classification losses $p(T|\hat{I}, \hat{C})$ and $p(Y|T, \hat{C}, \hat{A})$ on the CEVAE loss. However, since both latent confounders C and latent post-treatment variables M are correlated with both T and Y , these methods cannot disentangle C from M . An exception is NICE [38], which uses invariant risk minimization (IRM) [3] to find all causal parents of the outcome Y as the confounders, which makes it more robust in the LatentCorrelator case. However, since mediators M are also the causal parent of Y , the performance degrades substantially on the LatentMediator dataset. Although AFS [41] considers the existence of post-treatment variables M in the proxy X , it assumes that they can be separated from other variables in X in the observational space, and no relationship exists between the post-treatment variables and the outcome, so it still has poor performance in our setting since both assumptions are violated.

5.3 Sensitivity Analysis

In this part, we vary the number of confounders and post-treatment variables that generate proxy X in the Company (Age) and Company (Gender) datasets and compare CiVAE with the baseline TEDVAE in Fig. 3. Fig. 3 shows that the error is consistently lower for CiVAE. In addition, the error is comparatively higher when the number of confounders is low since the misidentification of latent post-treatment variables as confounders can have a comparatively larger influence on the ATE estimation. In addition, when the number of confounders becomes larger, the performance gap between CiVAE and TEDVAE gracefully shrinks.

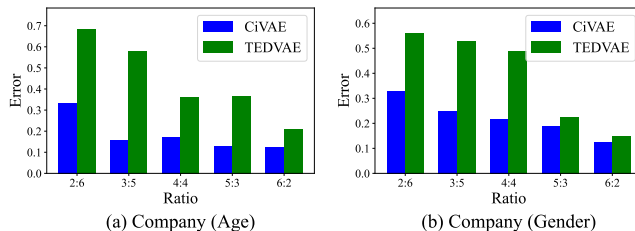


Figure 3: Error with different ratio of latent confounders and latent post-treatment variable in the latent space.

6 Conclusions

In this paper, we systematically investigate the latent post-treatment bias in causal inference from observational data. We first prove that unresolved latent post-treatment variables scrambled in the proxy of confounders can arbitrarily bias the ATE estimation. To address the bias, we proposed the Confounder-identifiable VAE (CiVAE), which, utilizing a mild assumption regarding the prior of latent factors, guarantees the identifiability of latent confounders up to bijective transformations. Finally, we show that controlling the latent confounders inferred by CiVAE can provide an unbiased estimation of the ATE. Experiments on both simulated and real-world datasets demonstrate that CiVAE has superior robustness to latent post-treatment bias compared to state-of-the-art methods.

373 **References**

- 374 [1] A. Acharya, M. Blackwell, and M. Sen. Explaining causal findings without bias: Detecting and
375 assessing direct effects. *American Political Science Review*, 110(3):512–529, 2016.
- 376 [2] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky. Tensor decompositions for
377 learning latent variable models. *Journal of Machine Learning Research*, 15:2773–2832, 2014.
- 378 [3] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz. Invariant risk minimization. *arXiv*
379 *preprint arXiv:1907.02893*, 2019.
- 380 [4] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians.
381 *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- 382 [5] L. Cheng, R. Guo, and H. Liu. Causal mediation analysis with hidden confounders. In
383 *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*,
384 pages 113–122, 2022.
- 385 [6] T. D. Cook, D. T. Campbell, and W. Shadish. *Experimental and quasi-experimental designs for*
386 *generalized causal inference*. Houghton Mifflin Boston, MA, 2002.
- 387 [7] P. Ding and L. W. Miratrix. To adjust or not to adjust? sensitivity analysis of m-bias and
388 butterfly-bias. *Journal of Causal Inference*, 3(1):41–57, 2015.
- 389 [8] J. K. Edwards, S. R. Cole, and D. Westreich. All your data are always missing: incorporating
390 bias due to measurement error into the potential outcomes framework. *International Journal of*
391 *Epidemiology*, 44(4):1452–1459, 2015.
- 392 [9] F. Elwert and C. Winship. Endogenous selection bias: The problem of conditioning on a collider
393 variable. *Annual review of sociology*, 40:31–53, 2014.
- 394 [10] T. A. Glass, S. N. Goodman, M. A. Hernán, and J. M. Samet. Causal inference in public health.
395 *Annual Review of Public Health*, 34:61–75, 2013.
- 396 [11] N. Hassanpour and R. Greiner. Counterfactual regression with importance sampling weights.
397 In *IJCAI*, pages 5880–5887, 2019.
- 398 [12] N. Hassanpour and R. Greiner. Learning disentangled representations for counterfactual
399 regression. In *International Conference on Learning Representations*, 2020.
- 400 [13] A. Hyvärinen and P. Dayan. Estimation of non-normalized statistical models by score matching.
401 *Journal of Machine Learning Research*, 6(4), 2005.
- 402 [14] K. Imai, L. Keele, and D. Tingley. A general approach to causal mediation analysis. *Psycholog-*
403 *ical Methods*, 15(4):309, 2010.
- 404 [15] G. W. Imbens and D. B. Rubin. *Causal inference in statistics, social, and biomedical sciences*.
405 Cambridge University Press, 2015.
- 406 [16] K. Jager, C. Zoccali, A. Macleod, and F. Dekker. Confounding: what it is and how to deal with
407 it. *Kidney international*, 73(3):256–260, 2008.
- 408 [17] F. Johansson, U. Shalit, and D. Sontag. Learning representations for counterfactual inference.
409 In *International Conference on Machine Learning*, pages 3020–3029, 2016.
- 410 [18] M. Kalisch and P. Bühlman. Estimating high-dimensional directed acyclic graphs with the
411 pc-algorithm. *Journal of Machine Learning Research*, 8(3), 2007.
- 412 [19] I. Khemakhem, D. Kingma, R. Monti, and A. Hyvarinen. Variational autoencoders and nonlinear
413 ICA: A unifying framework. In *International Conference on Artificial Intelligence and Statistics*,
414 pages 2207–2217. PMLR, 2020.
- 415 [20] G. King. A hard unsolved problem? post-treatment bias in big social science questions. In
416 *Hard Problems in Social Science” Symposium, April*, volume 10, 2010.

- 417 [21] G. King and L. Zeng. The dangers of extreme counterfactuals. *Political Analysis*, 14(2):131–159,
418 2006.
- 419 [22] D. Koller and N. Friedman. *Probabilistic graphical models: principles and techniques*. MIT
420 press, 2009.
- 421 [23] M. Kuroki and J. Pearl. Measurement bias and effect restoration in causal inference. *Biometrika*,
422 101(2):423–437, 2014.
- 423 [24] F. Li, K. L. Morgan, and A. M. Zaslavsky. Balancing covariates via propensity score weighting.
424 *Journal of the American Statistical Association*, 113(521):390–400, 2018.
- 425 [25] C. Louizos, U. Shalit, J. M. Mooij, D. Sontag, R. Zemel, and M. Welling. Causal effect inference
426 with deep latent-variable models. *Advances in Neural Information Processing Systems*, 30,
427 2017.
- 428 [26] C. Lu, Y. Wu, J. M. Hernández-Lobato, and B. Schölkopf. Invariant causal representation
429 learning for out-of-distribution generalization. In *International Conference on Learning Repre-*
430 *sentations*, 2021.
- 431 [27] Y. Lu and J. Lu. A universal approximation theorem of deep neural networks for expressing
432 probability distributions. In *Advances in Neural Information Processing Systems*, pages 3094–
433 3105, 2020.
- 434 [28] D. Madras, E. Creager, T. Pitassi, and R. Zemel. Fairness through causal awareness: Learning
435 causal latent-variable models for biased data. In *Proceedings of the Conference on Fairness,*
436 *Accountability, and Transparency*, pages 349–358, 2019.
- 437 [29] W. Miao, Z. Geng, and E. J. Tchetgen Tchetgen. Identifying causal effects with proxy variables
438 of an unmeasured confounder. *Biometrika*, 105(4):987–993, 2018.
- 439 [30] J. M. Montgomery, B. Nyhan, and M. Torres. How conditioning on posttreatment variables
440 can ruin your experiment and what to do about it. *American Journal of Political Science*,
441 62(3):760–775, 2018.
- 442 [31] J. Pearl. On measurement bias in causal inference. *arXiv preprint arXiv:1203.3504*, 2012.
- 443 [32] J. Pearl. Conditioning on post-treatment variables. *Journal of Causal Inference*, 3(1):131–137,
444 2015.
- 445 [33] S. J. Pocock, S. E. Assmann, L. E. Enos, and L. E. Kasten. Subgroup analysis, covariate
446 adjustment and baseline comparisons in clinical trial reporting: current practice and problems.
447 *Statistics in Medicine*, 21(19):2917–2930, 2002.
- 448 [34] M. Prosperri, Y. Guo, M. Sperrin, J. S. Koopman, J. S. Min, X. He, S. Rich, M. Wang, I. E.
449 Buchan, and J. Bian. Causal inference and counterfactual prediction in machine learning for
450 actionable healthcare. *Nature Machine Intelligence*, 2(7):369–375, 2020.
- 451 [35] K. J. Rothman, S. Greenland, T. L. Lash, et al. *Modern epidemiology*, volume 3. Wolters
452 Kluwer Health/Lippincott Williams & Wilkins Philadelphia, 2008.
- 453 [36] U. Shalit, F. D. Johansson, and D. Sontag. Estimating individual treatment effect: generalization
454 bounds and algorithms. In *International Conference on Machine Learning*, pages 3076–3085,
455 2017.
- 456 [37] C. Shi, D. Blei, and V. Veitch. Adapting neural networks for the estimation of treatment effects.
457 In *Advances in Neural Information Processing Systems*, 2019.
- 458 [38] C. Shi, V. Veitch, and D. M. Blei. Invariant representation learning for treatment effect estimation.
459 In *Uncertainty in Artificial Intelligence*, pages 1546–1555. PMLR, 2021.
- 460 [39] P. Spirtes, C. N. Glymour, R. Scheines, and D. Heckerman. *Causation, prediction, and search*.
461 MIT press, 2000.

- 462 [40] S. Wager and S. Athey. Estimation and inference of heterogeneous treatment effects using
463 random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.
- 464 [41] H. Wang, K. Kuang, H. Chi, L. Yang, M. Geng, W. Huang, and W. Yang. Treatment effect
465 estimation with adjustment feature selection. In *Proceedings of the 29th ACM SIGKDD*
466 *Conference on Knowledge Discovery and Data Mining*, pages 2290–2301, 2023.
- 467 [42] L. Yao, S. Li, Y. Li, M. Huai, J. Gao, and A. Zhang. Representation learning for treatment effect
468 estimation from observational data. In *Advances in Neural Information Processing Systems*,
469 volume 31, 2018.
- 470 [43] K. Zhang, J. Peters, D. Janzing, and B. Schölkopf. Kernel-based conditional independence test
471 and application in causal discovery. *arXiv preprint arXiv:1202.3775*, 2012.
- 472 [44] W. Zhang, L. Liu, and J. Li. Treatment effect estimation with disentangled latent factors. In
473 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 10923–10930,
474 2021.

475 Appendix

476 A Broader Impact

477 The proposed CiVAE is a universal model for causal effect estimation with observational data.
478 Although we use the Company job data that estimate the causal effects of *online working mode* to
479 *applicant statistics* as a real-world example, proxy-of-confounder-based methods have been heavily
480 used in other observational studies, which may be susceptible to latent post-treatment bias. Therefore,
481 we speculate that the proposed CiVAE will have a broader impact on causal inference community.

482 B Related Work

483 B.1 Post-Treatment Bias in Causal Inference

484 Bias due to accidentally controlling post-treatment variables, i.e., *post-treatment bias*, has long been
485 recognized as dangerous in causal effect estimation [20]. Back at 2005, Pearl [32] cautioned that
486 controlling more is not better, and uses the collider bias [9] and M-Bias [7] as two examples to
487 show that bias can be increased when controlling the post-treatment variables. Furthermore, [30]
488 show that indirect correlations between post-treatment variable M and outcome Y can still cause
489 bias. Recent works prove that even if M has no causal relationship with Y , controlling it can still
490 increase the variance of estimand [12]. However, most of these works study the post-treatment bias
491 in the observational space, where latent post-treatment variables that are mixed with confounders
492 to generate the observed covariates can be easily ignored by the researcher. Therefore, it motivates us to
493 develop CiVAE, which is robust to the latent post-treatment bias under mild assumptions.

494 B.2 Covariate Disentanglement

495 Recently, researchers have realized that directly controlling proxy of confounders \mathbf{X} may not be
496 safe, as variables other than confounders could lurk in the proxy and ruin the ATE estimation [12].
497 Traditional methods assume that the variables that generate \mathbf{X} are a mixture of confounders, adjusters,
498 and influencers [36], where adjusters should not be controlled as it can increase the estimation
499 variance [11]. Most methods rely on the fact that adjusters are correlated only with the treatment
500 to separate them from other variables [12, 44] (see Fig. (1)). This can also be used to remove post-
501 treatment variables that are not correlated with the outcome, which have similar statistics properties
502 with adjusters [41]. Here, a different work is NICE [38], which uses the fact that confounders and
503 influencers are direct causal parents of the outcome to find these variables with invariant learning as
504 the control set [3]. However, since mediators are also direct parents of the outcome, NICE is still not
505 robust to general post-treatment bias. Given that all above methods cannot satisfactorily address the
506 latent post-treatment in general cases, it is imperative to design the CiVAE, where confounders can
507 be identified and distinguished with latent post-treatment variables for unbiased adjustment.

508 C Theoretical Analysis

509 C.1 Proof of Lemma 3.1.

510 *Proof.* Let $\mathbf{Z} = f(\mathbf{X})$ and $z = f(\mathbf{x})$. If f is injective and differentiable *a.e.*, and f^\dagger is the
511 left-inverse, we have:

$$f_{Y|f(\mathbf{X})}(y|f(\mathbf{x})) = f_{Y|\mathbf{Z}}(y|z) = \frac{f_{Y,\mathbf{Z}}(y, z)}{f_{\mathbf{Z}}(z)} = \frac{f_{Y,\mathbf{X}}(y, f^\dagger(z))|\mathbf{J}_{f^\dagger}(z)|}{f_{\mathbf{X}}(f^\dagger(z))|\mathbf{J}_{f^\dagger}(z)|} = \frac{f_{Y,\mathbf{X}}(y, \mathbf{x})}{f_{\mathbf{X}}(\mathbf{x})} = f_{Y|\mathbf{X}}(y|\mathbf{x}), \quad (12)$$

512 where f and $f_{\cdot|\cdot}$ represent the marginal and conditional density function, respectively, and $\mathbf{J}_{f^\dagger}(z)$ is
513 the Jacobian matrix of function f^\dagger evaluated at z . Based on Eq. (12), we have:

$$\mathbb{E}[Y|\mathbf{X}] = \int \mathbf{y} \cdot f_{Y|\mathbf{X}}(\mathbf{y}|\mathbf{x}) d\mathbf{y} = \int \mathbf{y} \cdot f_{Y|\mathbf{Z}}(\mathbf{y}|z) d\mathbf{y} = \mathbb{E}[Y|\mathbf{Z} = z] = \mathbb{E}[Y|f(\mathbf{X}) = f(\mathbf{x})]. \quad (13)$$

514 □

515 **C.2 Proof of Corollary 3.3.**

516 *Proof.* For $\mathbf{X} = \mathbf{x}$, let $[c|\mathbf{m}] \doteq [f_C^\dagger(\mathbf{x})|f_M^\dagger(\mathbf{x})] \doteq f^\dagger(\mathbf{x}) = \mathbf{A}^\dagger(\mathbf{x} - \alpha_X)$, where \mathbf{A}^\dagger is the left
 517 inverse of the full column-rank matrix \mathbf{A} in Eq. (2), we have:

$$\begin{aligned}
 CATE(\mathbf{x}) &= \mathbb{E}[Y|T = 1, \mathbf{C} = f_C^\dagger(\mathbf{x})] - \mathbb{E}[Y|T = 0, \mathbf{C} = f_C^\dagger(\mathbf{x})] \\
 &= \mathbb{E}[Y|T = 1, \mathbf{C} = \mathbf{c}] - \mathbb{E}[Y|T = 0, \mathbf{C} = \mathbf{c}] \\
 &= \mathbb{E}[\alpha_Y + \tau \cdot T + \sum \theta_j \cdot M_j + \sum \kappa_i \cdot C_i | T = 1, \mathbf{C} = \mathbf{c}] \\
 &\quad - \mathbb{E}[\alpha_Y + \tau \cdot T + \sum \theta_j \cdot M_j + \sum \kappa_i \cdot C_i | T = 0, \mathbf{C} = \mathbf{c}] \\
 &= \alpha_Y + \tau \cdot \mathbb{E}[T|T = 1, \mathbf{C} = \mathbf{c}] + \sum \theta_j \cdot \mathbb{E}[M_j|T = 1, \mathbf{C} = \mathbf{c}] + \sum \kappa_i \cdot \mathbb{E}[C_i|T = 1, \mathbf{C} = \mathbf{c}] \\
 &\quad - \alpha_Y + \tau \cdot \mathbb{E}[T|T = 0, \mathbf{C} = \mathbf{c}] + \sum \theta_j \cdot \mathbb{E}[M_j|T = 0, \mathbf{C} = \mathbf{c}] + \sum \kappa_i \cdot \mathbb{E}[C_i|T = 0, \mathbf{C} = \mathbf{c}] \\
 &= \tau \cdot (1 - 0) + \sum \theta_j \cdot (\gamma_j \cdot (1 - 0)) + \sum \kappa_i \cdot (c_i - c_i) \\
 &= \tau + \sum \theta_j \cdot \gamma_j = \mathbb{E}[\tau + \sum \theta_j \cdot \gamma_j] = ATE,
 \end{aligned} \tag{14}$$

518 where the first equality is due to the definition of CATE in Eq. (2). In addition, the causal estimand
 519 and bias of a proxy-of-confounder-based causal inference model that controls the latent variable \mathbf{Z}
 520 inferred via $\mathbf{Z} = f(\mathbf{X}) = \mathbf{B}^T \mathbf{X}$ (where \mathbf{B} is also a full column-rank matrix) can be formulated as:

$$\begin{aligned}
 DCEV(\mathbf{B}^T \mathbf{x}) &= \mathbb{E}[Y|T = 1, \mathbf{Z} = \mathbf{B}^T \mathbf{x}] - \mathbb{E}[Y|T = 0, \mathbf{Z} = \mathbf{B}^T \mathbf{x}] \\
 &= \mathbb{E}[Y|T = 1, \mathbf{Z} = \mathbf{B}^T \alpha_X + \mathbf{B}^T \mathbf{A}[c|\mathbf{m}]] - \mathbb{E}[Y|T = 0, \mathbf{Z} = \mathbf{B}^T \alpha_X + \mathbf{B}^T \mathbf{A}[c|\mathbf{m}]] \\
 &\stackrel{(a)}{=} \mathbb{E}[Y|T = 1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] - \mathbb{E}[Y|T = 0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\
 &= \alpha_Y + \tau \cdot 1 + \sum \theta_j \cdot \mathbb{E}[M_j|T = 1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] + \sum \kappa_i \cdot \mathbb{E}[C_i|T = 1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\
 &\quad - \alpha_Y + \tau \cdot 0 + \sum \theta_j \cdot \mathbb{E}[M_j|T = 0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] + \sum \kappa_i \cdot \mathbb{E}[C_i|T = 0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\
 &= \tau \cdot (1 - 0) + \sum \theta_j \cdot (m_j - m_j) + \sum \kappa_i \cdot (c_i - c_i) \\
 &= \tau = \mathbb{E}[\tau] = \mathbb{E}[DCEV(\mathbf{B}^T \mathbf{X})],
 \end{aligned} \tag{15}$$

521 where step (a) is due to the fact that, since both \mathbf{A} and \mathbf{B} are full column-rank matrices, $\mathbf{B}^T \mathbf{A}$ is
 522 an invertible matrix, and the mapping $f = \mathbf{B}^T \alpha_X + \mathbf{B}^T \mathbf{A}$ is bijective. Therefore, we can invoke
 523 Lemma 3.1 and apply the left-inverse of f , i.e., $f^\dagger = (\mathbf{B}^T \mathbf{A})^{-1} - \mathbf{B}^T \alpha_X$, to the condition of the
 524 expectation. The rest steps are based on the structural causal equations defined in Eq. (2). \square

525 **C.3 Another Case of Linear SCM with Latent Correlators**

526 **Corollary C.1.** *For another Linear Structural Causal Model defined as follows*

$$\begin{aligned}
 T &\leftarrow \mathbb{1}(\alpha_T + \sum \beta_i \cdot C_i > a) \\
 M_j &\leftarrow \alpha_M + \gamma_j \cdot T + \phi_j \cdot U_j \\
 \mathbf{X} &\leftarrow \alpha_X + \mathbf{A}[M|\mathbf{C}] \\
 Y &\leftarrow \alpha_Y + \tau \cdot T + \sum \theta_j \cdot U_j + \sum \kappa_i \cdot C_i,
 \end{aligned} \tag{16}$$

527 where $f = \mathbf{A} \in \mathbb{R}^{K_X \times (K_C + K_M)}$ is a full column-rank matrix, the CATE, ATE, and the bias of
 528 proxy-of-confounder-based causal inference model that controls the latent variable \mathbf{Z} inferred via
 529 $\mathbf{Z} = f(\mathbf{X}) = \mathbf{B}^T \mathbf{X}$ can be formulated as follows:

$$\begin{aligned}
 ATE &= CATE = \tau \\
 \mathbb{E}[DCEV(\mathbf{Z} = \mathbf{B}^T \mathbf{X})] &= DCEV(\mathbf{Z} = \mathbf{B}^T \mathbf{X}) = \tau - \sum \frac{\theta_j \cdot \gamma_j}{\phi_j} \\
 Bias &= ATE - \mathbb{E}[DCEV(\mathbf{B}^T \mathbf{X})] = \sum \frac{\theta_j \cdot \gamma_j}{\phi_j},
 \end{aligned} \tag{17}$$

530 where $\mathbf{B} \in \mathbb{R}^{K_X \times (K_C + K_M)}$ is another full column-rank matrix. Since $\sum \frac{\theta_j \cdot \gamma_j}{\phi_j}$ is arbitrary, the
 531 estimator $\mathbb{E}[DCEV(\mathbf{B}^T \mathbf{X})]$ is arbitrarily biased for the estimation of ATE.

532 *Proof.* The proof of the CATE and ATE is trivial. The causal estimand and the bias of a proxy-
 533 of-confounder-based causal inference model that controls the latent variables \mathbf{Z} inferred via $\mathbf{Z} =$
 534 $\tilde{f}(\mathbf{X}) = \mathbf{B}^T \mathbf{X}$ (where \mathbf{B} is also a full column-rank matrix) can be formulated as follows:

$$\begin{aligned}
 DCEV(\mathbf{B}^T \mathbf{x}) &= \mathbb{E}[Y|T = 1, \mathbf{Z} = \mathbf{B}^T \mathbf{x}] - \mathbb{E}[Y|T = 0, \mathbf{Z} = \mathbf{B}^T \mathbf{x}] \\
 &= \mathbb{E}[Y|T = 1, \mathbf{Z} = \boldsymbol{\alpha}_X + \mathbf{B}^T \mathbf{A}[\mathbf{c}|\mathbf{m}]] - \mathbb{E}[Y|T = 0, \mathbf{Z} = \boldsymbol{\alpha}_X + \mathbf{B}^T \mathbf{A}[\mathbf{c}|\mathbf{m}]] \\
 &\stackrel{(a)}{=} \mathbb{E}[Y|T = 1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] - \mathbb{E}[Y|T = 0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\
 &= \alpha_Y + \tau \cdot 1 + \sum \theta_j \cdot \mathbb{E}[U_j|T = 1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] + \sum \kappa_i \cdot \mathbb{E}[C_i|T = 1, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\
 &\quad - \alpha_Y + \tau \cdot 0 + \sum \theta_j \cdot \mathbb{E}[U_j|T = 0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] + \sum \kappa_i \cdot \mathbb{E}[C_i|T = 0, \mathbf{C} = \mathbf{c}, \mathbf{M} = \mathbf{m}] \\
 &= \tau \cdot (1 - 0) + \sum \theta_j \cdot (\phi_j^{-1} \cdot (m_j - \alpha_M - \gamma_j) - \phi_j^{-1} \cdot (m_j - \alpha_M)) + \sum \kappa_i \cdot (c_i - c_i) \\
 &= \tau - \sum \frac{\theta_j \cdot \gamma_j}{\phi_j} = \mathbb{E} \left[\tau - \sum \frac{\theta_j \cdot \gamma_j}{\phi_j} \right] = \mathbb{E}[DCEV(\mathbf{B}^T \mathbf{X})],
 \end{aligned}$$

(18)

535

□

536 where step (a) and the rest of the proof follow the same logic as the proof in Section 3.3.

537 C.4 Proof of Theorem 4.1

538 The strict definitions of the exponential family, strong exponential (which is assumed for the factorized
 539 part of the conditional prior), and identifiability follow [19, 26], and can be referred to in Appendix
 540 E, F of [26], which we omit to avoid redundancy. The proof of Theorem 4.1 is largely based on the
 541 NF-iVAE paper [26], where most of the details can be found, with the new assumption introduced in
 542 CiVAE that each $\mathbf{S}_{f,i}$ has at least one invertible dimension incorporated to ensure that each dimension
 543 of the inferred latent variables is a bijective transformation of the corresponding true latent variable.

544 C.4.1 PART I

545 **Step I.** In this step, we transform the equality of noisy conditional marginal distribution of \mathbf{X} given
 546 Y, T of two models with parameter $\boldsymbol{\theta}$, $\boldsymbol{\theta} \in \Theta$ into the equality of noise-free distributions.

$$\begin{aligned}
 p_{\boldsymbol{\theta}}(\mathbf{X} | Y, T) &= p_{\tilde{\boldsymbol{\theta}}}(\mathbf{X} | Y, T) \\
 \implies \int_{\mathbf{Z}} p_f(\mathbf{X} | \mathbf{Z}) p_{S, \lambda}(\mathbf{Z} | Y, T) d\mathbf{Z} &= \int_{\mathbf{Z}} p_{\tilde{f}}(\mathbf{X} | \mathbf{Z}) p_{\tilde{S}, \tilde{\lambda}}(\mathbf{Z} | Y, T) d\mathbf{Z} \\
 \implies \int_{\mathbf{Z}} p_{\epsilon}(\mathbf{X} - f(\mathbf{Z})) p_{S, \lambda}(\mathbf{Z} | Y, T) d\mathbf{Z} &= \int_{\mathbf{Z}} p_{\epsilon}(\mathbf{X} - \tilde{f}(\mathbf{Z})) p_{\tilde{S}, \tilde{\lambda}}(\mathbf{Z} | Y, T) d\mathbf{Z} \\
 \stackrel{(a)}{\implies} \int_{\mathcal{X}} p_{\epsilon}(\mathbf{X} - \overline{\mathbf{X}}) p_{S, \lambda}(f^{\dagger}(\overline{\mathbf{X}}) | Y, T) \text{vol}(\mathbf{J}_{f^{\dagger}}(\overline{\mathbf{X}})) d\overline{\mathbf{X}} &= \\
 \int_{\mathcal{X}} p_{\epsilon}(\mathbf{X} - \overline{\mathbf{X}}) p_{\tilde{S}, \tilde{\lambda}}(\tilde{f}^{\dagger}(\overline{\mathbf{X}}) | Y, T) \text{vol}(\mathbf{J}_{\tilde{f}^{\dagger}}(\overline{\mathbf{X}})) d\overline{\mathbf{X}} & \\
 \stackrel{(b)}{\implies} \int_{\mathbb{R}^d} p_{\epsilon}(\mathbf{X} - \overline{\mathbf{X}}) \tilde{p}_{f, S, \lambda, Y, T}(\overline{\mathbf{X}}) d\overline{\mathbf{X}} &= \int_{\mathbb{R}^d} p_{\epsilon}(\mathbf{X} - \overline{\mathbf{X}}) \tilde{p}_{\tilde{f}, \tilde{S}, \tilde{\lambda}, \tilde{Y}, \tilde{T}}(\overline{\mathbf{X}}) d\overline{\mathbf{X}} \\
 \implies (\tilde{p}_{f, S, \lambda, Y, T} * p_{\epsilon})(\mathbf{X}) &= (\tilde{p}_{\tilde{f}, \tilde{S}, \tilde{\lambda}, \tilde{Y}, \tilde{T}} * p_{\epsilon})(\mathbf{X}) \\
 \stackrel{(c)}{\implies} F[\tilde{p}_{f, S, \lambda, Y, T}](\boldsymbol{\omega}) \varphi_{\epsilon}(\boldsymbol{\omega}) &= F[\tilde{p}_{\tilde{f}, \tilde{S}, \tilde{\lambda}, \tilde{Y}, \tilde{T}}](\boldsymbol{\omega}) \varphi_{\epsilon}(\boldsymbol{\omega}) \\
 \stackrel{(d)}{\implies} F[\tilde{p}_{f, S, \lambda, Y, T}](\boldsymbol{\omega}) &= F[\tilde{p}_{\tilde{f}, \tilde{S}, \tilde{\lambda}, \tilde{Y}, \tilde{T}}](\boldsymbol{\omega}) \\
 \implies \tilde{p}_{f, S, \lambda, Y, T}(\mathbf{X}) &= \tilde{p}_{\tilde{f}, \tilde{S}, \tilde{\lambda}, \tilde{Y}, \tilde{T}}(\mathbf{X}).
 \end{aligned}$$

(19)

547 Step (a) is based on the rule of change-of-variable, where $\text{vol}(\mathbf{A}) = \sqrt{\det(\mathbf{A}^T \mathbf{A})}$. In step (b),
 548 we define $\tilde{p}_{f, \mathcal{S}, \lambda, Y, T}(\mathbf{X}) \triangleq p_{\mathcal{S}, \lambda}(f^\dagger(\mathbf{X}) | Y, T) \text{vol}(\mathbf{J}_{f^\dagger}(\mathbf{X})) \mathbb{I}_{\mathcal{X}}(\mathbf{X})$. In step (c), we use $F[\cdot]$ to
 549 denote the Fourier transform. In step (d), we drop $\varphi_\varepsilon(\omega)$ as it is non-zero *a.e.* (see Assumption 3).

550 **Step II.** In this step, we transform the equality of the noise-free distributions into the relationship of
 551 the sufficient statistics \mathbf{S} and $\tilde{\mathbf{S}}$. By taking logarithm of both sides of Eq. (19), we have:

$$\begin{aligned} & \log \text{vol}(J_{f^\dagger}(\mathbf{X})) + \log \mathcal{Q}(f^\dagger(\mathbf{X})) - \log \mathcal{C}(Y, T) + \langle \mathbf{S}(f^\dagger(\mathbf{X})), \lambda(Y, T) \rangle \\ &= \log \text{vol}(J_{\tilde{f}^\dagger}(\mathbf{X})) + \log \tilde{\mathcal{Q}}(\tilde{f}^\dagger(\mathbf{X})) - \log \tilde{\mathcal{C}}(Y, T) + \langle \tilde{\mathbf{S}}(\tilde{f}^\dagger(\mathbf{X})), \tilde{\lambda}(Y, T) \rangle. \end{aligned} \quad (20)$$

552 Let $(Y, T)_0, \dots, (Y, T)_k$ be the $k + 1$ distinct points defined in Assumption 3 - (iv). We obtain $k + 1$
 553 equations by evaluating the Eq. (20) at these points, where the first equation is subtracted from the
 554 remaining ones, which leads to the following equation system:

$$\begin{aligned} & \langle \mathbf{S}(f^\dagger(\mathbf{X})), \lambda((Y, T)_l) - \lambda((Y, T)_0) \rangle + \log \frac{\mathcal{C}((Y, T)_0)}{\mathcal{C}((Y, T)_l)} \\ &= \langle \tilde{\mathbf{S}}(\tilde{f}^\dagger(\mathbf{X})), \tilde{\lambda}((Y, T)_l) - \tilde{\lambda}((Y, T)_0) \rangle + \log \frac{\tilde{\mathcal{C}}((Y, T)_0)}{\tilde{\mathcal{C}}((Y, T)_l)}, \quad l = 1, \dots, k. \end{aligned} \quad (21)$$

555 Let \mathbf{L} be the invertible matrix defined in Assumption 3 - (iv) and $\tilde{\mathbf{L}}$ be the counterpart for $\tilde{\lambda}$, if we
 556 summarize all terms irrelevant to \mathbf{X} into a constant \mathbf{b} , we have:

$$\begin{aligned} & \mathbf{L}^T \mathbf{S}(f^\dagger(\mathbf{X})) = \tilde{\mathbf{L}}^T \tilde{\mathbf{S}}(\tilde{f}^\dagger(\mathbf{X})) + \mathbf{b} \\ & \implies \mathbf{S}(f^\dagger(\mathbf{X})) = \mathbf{A} \tilde{\mathbf{S}}(\tilde{f}^\dagger(\mathbf{X})) + \mathbf{c}, \end{aligned} \quad (22)$$

557 where $\mathbf{A} = \mathbf{L}^{-T} \tilde{\mathbf{L}} \in \mathbb{R}^{k \times k}$, and $\mathbf{c} = \mathbf{L}^{-T} \mathbf{b} \in \mathbb{R}^k$.

558 **Step III.** Ideally, to prove the element-wise bijective identifiability of the latent variables \mathbf{Z} , the
 559 transformation of the sufficient statistics \mathbf{S} derived in Eq. (22) should be bijective. We claim that if
 560 the conditional prior $p_{\mathcal{S}, \lambda}(\mathbf{Z} | Y, T)$ is strongly exponential and \mathbf{L} is invertible, $\tilde{\mathbf{L}}$ and \mathbf{A} must also
 561 be invertible. The proof is omitted, and can be referred to in Appendix H.1.1 of [26].

562 C.4.2 PART II

563 In this part, we prove that, if Assumptions 1, 2 and 3 hold, we can identify the factorized part
 564 of the sufficient statistics $\mathbf{S}(\mathbf{Z})$, i.e., $\mathbf{S}_f(\mathbf{Z})$, up to permutation and element-wise transformation.
 565 Specifically, if we use \mathbf{v} to denote the composite map $\tilde{f}^\dagger \circ f : \mathcal{Z} \rightarrow \mathcal{Z}$, Eq. (22) can be rewritten into:

$$\mathbf{S}(\mathbf{Z}) = \mathbf{A} \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z})) + \mathbf{c}. \quad (23)$$

566 We aim to prove that \mathbf{A} in Eq. (23) is a block permutation matrix.

567 **Step I.** We start by showing that \mathbf{v} is a component-wise function. If we differentiate both sides of Eq.
 568 (23) with respect to Z_s and Z_t , where $s \neq t$, we have:

$$\begin{aligned} \frac{\partial \mathbf{S}(\mathbf{Z})}{\partial Z_s} &= \mathbf{A} \sum_{i=1}^{K_Z} \frac{\partial \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_s} \\ \frac{\partial^2 \mathbf{S}(\mathbf{Z})}{\partial Z_s \partial Z_t} &= \mathbf{A} \sum_{i=1}^{K_Z} \sum_{j=1}^{K_Z} \frac{\partial^2 \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z}) \partial v_j(\mathbf{Z})} \cdot \frac{\partial v_j(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_s} + \mathbf{A} \sum_{i=1}^{K_Z} \frac{\partial \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})} \cdot \frac{\partial^2 v_i(\mathbf{Z})}{\partial Z_s \partial Z_t}. \end{aligned} \quad (24)$$

569 Note that for the factorized part of the sufficient statistics \mathbf{S} , i.e., \mathbf{S}_f , all *cross-derivatives* are zero,
 570 and for the non-factorized part of \mathbf{S} , i.e., \mathbf{S}_{nf} , which is a neural network with ReLU activation (i.e.,
 571 linear *a.e.*), all *second-order derivatives* are zero. Therefore, the *second order cross-derivatives* on
 572 the LHS. of Eq. (24) are zero, which leads to the following equality:

$$\mathbf{0} = \mathbf{A} \sum_{i=1}^{K_Z} \frac{\partial^2 \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})^2} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_i(\mathbf{Z})}{\partial Z_s} + \mathbf{A} \sum_{i=1}^{K_Z} \frac{\partial \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_i(\mathbf{Z})} \cdot \frac{\partial^2 v_i(\mathbf{Z})}{\partial Z_s \partial Z_t}. \quad (25)$$

573 Eq. (25) can be written into the matrix-vector product form as follows:

$$\mathbf{0} = \mathbf{A}\tilde{\mathbf{S}}''(\mathbf{Z})\mathbf{v}'_{s,t}(\mathbf{Z}) + \mathbf{A}\tilde{\mathbf{S}}'(\mathbf{Z})\mathbf{v}''_{s,t}(\mathbf{Z}), \quad (26)$$

where

$$\begin{aligned} \tilde{\mathbf{S}}''(\mathbf{Z}) &= \left[\frac{\partial^2 \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_1(\mathbf{Z})^2}, \dots, \frac{\partial^2 \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_{K_Z}(\mathbf{Z})^2} \right] \in \mathbb{R}^{k \times K_Z}, \\ \mathbf{v}'_{s,t}(\mathbf{Z}) &= \left[\frac{\partial v_1(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_1(\mathbf{Z})}{\partial Z_s}, \dots, \frac{\partial v_{K_Z}(\mathbf{Z})}{\partial Z_t} \cdot \frac{\partial v_{K_Z}(\mathbf{Z})}{\partial Z_s} \right]^T \in \mathbb{R}^{K_Z}, \end{aligned}$$

and

$$\begin{aligned} \tilde{\mathbf{S}}'(\mathbf{Z}) &= \left[\frac{\partial \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_1(\mathbf{Z})}, \dots, \frac{\partial \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))}{\partial v_{K_Z}(\mathbf{Z})} \right] \in \mathbb{R}^{k \times K_Z}, \\ \mathbf{v}''_{s,t}(\mathbf{Z}) &= \left[\frac{\partial^2 v_1(\mathbf{Z})}{\partial Z_s \partial Z_t}, \dots, \frac{\partial^2 v_{K_Z}(\mathbf{Z})}{\partial Z_s \partial Z_t} \right]^T \in \mathbb{R}^{K_Z}. \end{aligned}$$

574 If we denote the concatenation as $\tilde{\mathbf{S}}'''(\mathbf{Z}) = \left[\tilde{\mathbf{S}}''(\mathbf{Z}), \tilde{\mathbf{S}}'(\mathbf{Z}) \right] \in \mathbb{R}^{k \times 2K_Z}$ and $\mathbf{v}''_{s,t}(\mathbf{Z}) =$
 575 $\left[\mathbf{v}'_{s,t}(\mathbf{Z})^T, \mathbf{v}''_{s,t}(\mathbf{Z})^T \right]^T \in \mathbb{R}^{2K_Z}$, we have:

$$\mathbf{0} = \mathbf{A}\tilde{\mathbf{S}}'''(\mathbf{Z})\mathbf{v}''_{s,t}(\mathbf{Z}). \quad (27)$$

576 Finally, if we denote the rows of $\tilde{\mathbf{S}}'''(\mathbf{Z})$ that correspond to the factorized part of \mathbf{S} by $\tilde{\mathbf{S}}'_f(\mathbf{Z})$,
 577 according to Lemma 5 of the iVAE paper [19] and the assumption that $k \geq 2K_Z$, we have that the
 578 rank of $\tilde{\mathbf{S}}'_f(\mathbf{Z})$ is $2K_Z$. Since $k \geq 2K_Z$, the rank of $\tilde{\mathbf{S}}'''(\mathbf{Z})$ is also $2K_Z$. Since the rank of \mathbf{A} is k ,
 579 the rank of $\mathbf{A}\tilde{\mathbf{S}}'''(\mathbf{Z})$ is $2K_Z$, which implies that $\mathbf{v}''_{s,t}(\mathbf{Z}) \in \mathbb{R}^{2K_Z}$ is a zero vector. Therefore, we
 580 have $\mathbf{v}'_{s,t}(\mathbf{Z}) = \mathbf{0}, \forall s \neq t$, and we have demonstrated that \mathbf{v} is a component-wise function.

581 **Step II.** Based on **Step I**, we demonstrate that \mathbf{A} is a block permutation matrix. Without loss of gen-
 582 erality, we assume that the permutation in \mathbf{v} is Identity, where $\mathbf{v}(\mathbf{Z}) = [v_1(Z_1), \dots, v_{K_Z}(Z_{K_Z})]^T$
 583 and each v_i is a nonlinear univariate scalar function. Since f and \tilde{f} are injective, \mathbf{v} is bijective and
 584 $\mathbf{v}^{-1}(\mathbf{Z}) = [v_1^{-1}(Z_1), \dots, v_{K_Z}^{-1}(Z_{K_Z})]^T$. If we denote $\bar{\mathbf{S}}(\mathbf{v}(\mathbf{Z})) = \tilde{\mathbf{S}}(\mathbf{v}(\mathbf{Z})) + \mathbf{A}^{-1}\mathbf{c}$, Eq. (23)
 585 can be reformulated as $\mathbf{S}(\mathbf{Z}) = \mathbf{A}\bar{\mathbf{S}}(\mathbf{v}(\mathbf{Z}))$. We then apply \mathbf{v}^{-1} to \mathbf{Z} on both sides, which gives

$$\mathbf{S}(\mathbf{v}^{-1}(\mathbf{Z})) = \mathbf{A}\bar{\mathbf{S}}(\mathbf{Z}). \quad (28)$$

586 Let t be the index of an entry in \mathbf{S} that corresponds to the factorized part \mathbf{S}_f . For all $s \neq t$, we have:

$$0 = \frac{\partial \mathbf{S}(\mathbf{v}^{-1}(\mathbf{Z}))_t}{\partial Z_s} = \sum_{j=1}^k a_{tj} \frac{\partial \bar{\mathbf{S}}(\mathbf{Z})_j}{\partial Z_s}. \quad (29)$$

587 Since the entries of $\bar{\mathbf{S}}$ are linearly independent, a_{tj} is zero for any j such that $\frac{\partial \bar{\mathbf{S}}(\mathbf{Z})_j}{\partial Z_s} \neq 0$. This
 588 includes the entries S_j that correspond to (1) the factorized part that does not depend on Z_t ; and (2)
 589 the non-factorized part \mathbf{S}_{n_f} . Therefore, when t is the index of an entry in the sufficient statistics \mathbf{S}
 590 that corresponds to factor i in the factorized part \mathbf{S}_f , i.e., $\mathbf{S}_{f,i}$, the only non-zero a_{tj} are the ones that
 591 map between $\mathbf{S}_{f,i}(Z_i)$ and $\bar{\mathbf{S}}_{f,i}(v_i(Z_i))$. Therefore, we can construct an invertible submatrix \mathbf{A}'_i
 592 with all non-zero elements a_{tj} for all t that corresponds to factor i , such that

$$\mathbf{S}_{f,i}(Z_i) = \mathbf{A}'_i \bar{\mathbf{S}}_{f,i}(v_i(Z_i)) = \mathbf{A}'_i \tilde{\mathbf{S}}_{f,i}(v_i(Z_i)) + \mathbf{c}_i, \quad i = 1, \dots, K_Z, \quad (30)$$

593 where \mathbf{c}_i denotes the corresponding elements of \mathbf{c} . Eq. (30) means that for each $i = 1, \dots, K_Z$,
 594 the matrix block \mathbf{A}'_i of \mathbf{A} affinely transforms the i -specific sufficient statistics vector $\mathbf{S}_{f,i}(Z_i)$ into
 595 $\tilde{\mathbf{S}}_{f,i}(v_i(Z_i))$. In addition, there is also an additional block \mathbf{A}' that affinely transforms $\mathbf{S}_{n_f}(\mathbf{Z})$ in
 596 into $\mathbf{S}_{n_f}(\mathbf{v}(\mathbf{Z}))$. This completes the proof that \mathbf{A} is a block permutation matrix.

597 C.4.3 PART III

598 Let $\tilde{Z}_i = v_i(Z_i) = \tilde{f}^\dagger(\mathbf{X})_i$ be the i th inferred latent variable. Assume again that the permutation in
 599 \mathbf{v} is Identity. In this part, we prove that if Assumption 2 holds, each inferred latent variable \tilde{Z}_i is the
 600 bijective transformation of the true latent variable. The proof is as follows.

601 *Proof.* Plugging \tilde{Z}_i into Eq. (30), we have:

$$\mathbf{S}_{f,i}(Z_i) = \mathbf{A}'_i \tilde{\mathbf{S}}_{f,i}(\tilde{Z}_i). \quad (31)$$

602 According to Assumption 2, there exists one dimension of $\mathbf{S}_{f,i}$, i.e., j , such that $S_{f,i,j}$ is bijective.
 603 This implies that $\mathbf{S}_{f,i}$ is injective, and therefore it has a left-inverse $\mathbf{S}_{f,i}^\dagger$. we apply $\mathbf{S}_{f,i}^\dagger$ to both sides
 604 of Eq. (31), which gives:

$$Z_i = \mathbf{S}_{f,i}^\dagger \mathbf{A}'_i \tilde{\mathbf{S}}_{f,i}(\tilde{Z}_i). \quad (32)$$

605 Since \mathbf{A}'_i is a block of an invertible block permutation matrix, \mathbf{A}_i is also an invertible matrix, and
 606 therefore \mathbf{A}'_i is a bijective mapping. In addition, since $\tilde{\mathbf{S}}_{f,i}$ is injective, $\tilde{\mathbf{S}}_{f,i}$ is also injective, and
 607 therefore the composite map $\mathbf{S}_{f,i}^\dagger \mathbf{A}'_i \tilde{\mathbf{S}}_{f,i} : \mathbb{R} \rightarrow \mathbb{R}$ that applies on \tilde{Z}_i is a bijective. This completes
 608 the proof that each inferred latent variable \tilde{Z}_i is the bijective transformation of the true latent variable
 609 in the case of no noise, where $\mathbf{Z} = f^\dagger(\mathbf{X})$ are the true latent variables. If noise ε exists, the posterior
 610 distribution of the latent variables can be identified up to an analogous bijective indeterminacy. \square

611 C.4.4 Consistency

612 *Proof.* If the family of the variational posterior $q_\phi(\mathbf{Z}|\mathbf{X}, Y, T)$ contains the true posterior
 613 $p_\theta(\mathbf{Z}|\mathbf{X}, Y, T)$, then by optimizing the loss of Eq. (9) (with the KL term replaced by the score match-
 614 ing loss defined in Eq. (10)) over its parameter ϕ , the score matching term will eventually vanish.
 615 Therefore, the ELBO term in Eq. (9) will be equal to the log-likelihood. Under this circumstance,
 616 CiVAE inherits all the properties of maximum likelihood estimation (MLE). Since the identifiability
 617 of CiVAE is guaranteed up to permutation and component-wise bijective transformation of the latent
 618 variables, the consistency property of MLE means that the model will converge to the true parameter
 619 θ^* up to such mild indeterminacy of the latent variables in the limit of infinite data. \square

620 C.5 Proof of Theorem 4.2

621 *Proof.* Let \mathbf{C} be the true latent confounders and $\tilde{\mathbf{C}}$ be the transformed confounders, where the
 622 transformation function f is bijective and differentiable *a.e.* Let f^{-1} denote its inverse. The ATE
 623 estimator that controls transformed confounders $\tilde{\mathbf{C}}$ can be formulated as:

$$DEV(\tilde{\mathbf{C}}) = \mathbb{E}_{p(\tilde{\mathbf{C}})}[\mathbb{E}[Y|T = 1, \tilde{\mathbf{C}} = \tilde{\mathbf{c}}] - \mathbb{E}[Y|T = 0, \tilde{\mathbf{C}} = \tilde{\mathbf{c}}]]. \quad (33)$$

624 Specifically, for the continuous case where density functions exist, for each term, we have:

$$\mathbb{E}_{p(\tilde{\mathbf{C}})}[\mathbb{E}[Y|T = t, \tilde{\mathbf{C}} = \tilde{\mathbf{c}}]] = \int f_{\tilde{\mathbf{C}}}(\tilde{\mathbf{c}}) \int y \cdot f_{Y|T, \tilde{\mathbf{C}}}(y|t, \tilde{\mathbf{c}}) dy d\tilde{\mathbf{c}}. \quad (34)$$

625 For the marginal density $f_{\tilde{\mathbf{C}}}(\tilde{\mathbf{c}})$, the following equality holds:

$$f_{\tilde{\mathbf{C}}}(\tilde{\mathbf{c}}) = f_{\mathbf{C}}(f^{-1}(\tilde{\mathbf{c}})) |J_{f^{-1}}(\tilde{\mathbf{c}})| = f_{\mathbf{C}}(\mathbf{c}) |J_{f^{-1}}(\tilde{\mathbf{c}})|. \quad (35)$$

626 As for the conditional density $f_{Y|T, \tilde{\mathbf{C}}}(y|t, \tilde{\mathbf{c}})$, since f is bijective, according to Eq. (12), we have:

$$f_{Y|T, \tilde{\mathbf{C}}}(y|t, \tilde{\mathbf{c}}) = f_{Y|T, \mathbf{C}}(y|t, \mathbf{c}). \quad (36)$$

627 Combining Eqs. (35) and (36), and given that $d\tilde{\mathbf{c}} = |J_f(\mathbf{c})|d\mathbf{c}$, we have:

$$\begin{aligned} (34) &= \int f_{\mathbf{C}}(\mathbf{c}) |\mathbf{J}_{f^{-1}}(\tilde{\mathbf{c}})| \int y \cdot f_{Y|T, \mathbf{C}}(y|t, \mathbf{c}) dy |\mathbf{J}_f(\mathbf{c})| d\mathbf{c} \\ &= |\mathbf{J}_{f^{-1}}(\tilde{\mathbf{c}})| \cdot |\mathbf{J}_f(\mathbf{c})| \int f_{\mathbf{C}}(\mathbf{c}) \int y \cdot f_{Y|T, \mathbf{C}}(y|t, \mathbf{c}) dy d\mathbf{c} \\ &\stackrel{(a)}{=} \int f_{\mathbf{C}}(\mathbf{c}) \int y \cdot f_{Y|T, \mathbf{C}}(y|t, \mathbf{c}) dy d\mathbf{c} \\ &= \mathbb{E}_{p(\mathbf{C})}[\mathbb{E}[Y|T = t, \mathbf{C} = \mathbf{c}]], \end{aligned} \quad (37)$$

Table 2: Comparison of CiVAE with baselines when intra-interactions among M exist.

Dataset	LatentMediator		LatentCorrelator		Company (Age)		Company (Gender)	
	ATE.	Err.	ATE.	Err.	ATE.	Err.	ATE.	Err.
CEVAE	1.627 ± 0.549	2.627	2.659 ± 0.302	1.353	0.152 ± 0.027	0.420	-0.225 ± 0.044	-0.144
TEDVAE	1.653 ± 0.511	2.042	2.827 ± 0.259	1.521	0.180 ± 0.047	0.448	-0.189 ± 0.012	-0.108
CiVAE	-0.350 ± 0.695	1.785	1.785 ± 0.481	0.479	-0.073 ± 0.101	0.195	-0.136 ± 0.087	-0.055
True ATE	-1.000 ± 0.000	0.000	1.306 ± 0.000	0.000	-0.268 ± 0.000	0.000	-0.081 ± 0.000	0.000

Table 3: Comparison of CiVAE with baselines when inter-interactions between C and M exist.

Dataset	LatentMediator		LatentCorrelator		Company (Age)		Company (Gender)	
	ATE.	Err.	ATE.	Err.	ATE.	Err.	ATE.	Err.
CEVAE	2.070 ± 0.279	3.070	2.831 ± 0.398	1.831	0.094 ± 0.061	0.362	-0.192 ± 0.015	-0.111
TEDVAE	1.743 ± 0.307	2.743	2.954 ± 0.763	1.954	0.109 ± 0.116	0.377	-0.212 ± 0.019	-0.131
CiVAE	-0.716 ± 0.523	0.284	1.385 ± 0.660	0.385	-0.041 ± 0.144	0.227	-0.129 ± 0.064	-0.048
True ATE	-1.000 ± 0.000	0.000	1.000 ± 0.000	0.000	-0.268 ± 0.000	0.000	-0.081 ± 0.000	0.000

628 where the term $|J_{f^{-1}}(\tilde{c})| \cdot |J_f(c)|$ vanishes in step (a) as the two factors have the product of one.
 629 Therefore, if we plug Eq. (37) into Eq. (33), it leads to the following equality:

$$\begin{aligned}
 DEV(\tilde{C}) &= \mathbb{E}_{p(\tilde{C})}[\mathbb{E}[Y|T=1, \tilde{C}=\tilde{c}] - \mathbb{E}[Y|T=0, \tilde{C}=\tilde{c}]] \\
 &= \mathbb{E}_{p(C)}[\mathbb{E}[Y|T=1, C=c] - \mathbb{E}[Y|T=0, C=c]] = DEV(C) = ATE,
 \end{aligned}
 \tag{38}$$

630 where the last step is due to Eq. (2) in Definition 2, which completes our proof that controlling
 631 bijectively transformed confounders provides an unbiased estimation of ATE. \square

632 D Extending CiVAE to address Latent Interactions

633 In this section, we extend CiVAE to more general cases where interactions exist among the latent
 634 confounders C and the latent post-treatment variables M . Here, we note that the identification
 635 of latent confounders C in CiVAE is achieved in two steps. (i) CiVAE *individually* identifies
 636 latent variables $[C, M]$ that generate X in inferred Z (but which dims of Z correspond to C
 637 or M is unknown). (ii) pairwise independence test to identify C . Since Assumption 2 allows
 638 arbitrary dependence among C and M , step (i) still holds when interactions among $[C, M]$ exist.
 639 To distinguish C in these cases, we can use more general causal discovery algorithms, e.g., the
 640 PC algorithm [18] in the second step. In this section, we consider two cases of interaction: (i)
 641 Intra-Interaction among mediators, and (ii) Inter-Interaction among mediators and confounders.

642 D.1 Intra-Interactions among Latent Mediators

643 In this subsection, we discuss the case where latent post-treatment variables M interact with each
 644 other. Since in this case, M cannot causally influence the latent confounders C (otherwise C will be
 645 post-treatment), and the PC algorithm orients edges in causal graphs via colliders, latent confounders
 646 can still be identified from the inferred Z as they form colliders with the treatment T .

647 To empirically verify the claim, we extend the simulated datasets described in Section 5.1, where we
 648 make (i) T directly affects M_1 , (ii) M_1 affects M_2 , and (iii) M_1, M_2 affect M_3 . The coefficients are
 649 randomly sampled from $\mathcal{N}(0, 1/3)$. In step (ii), we use the PC algorithm [18] to identify C from
 650 the inferred Z . The results in Table 2 demonstrate that the adapted CiVAE is still significantly more
 651 robust to latent post-treatment bias compared to CEVAE and TEDVAE, which empirically verify our
 652 claim that PC-adapted CiVAE can address the interaction among post-treatment variables.

653 D.2 Inter-Interactions between Latent Mediators and Latent Confounders

654 In this subsection, we discuss another case where inter-interactions exist between latent confounders
 655 C and latent post-treatment variables M . Since in this case, M still cannot causally influence C
 656 (otherwise C will be post-treatment), and the PC algorithm orients edges in causal graph via colliders,
 657 latent confounders C can still be identified from Z as they form colliders with the treatment T .

658 To verify the claim, we extend the simulated datasets described in Section 5.1 to allow each latent
659 confounder $C_i \in \mathbb{R}^3$ to determine $M \in \mathbb{R}^3$. The coefficients are randomly sampled from $\mathcal{N}(0, 1/3)$.
660 In step **(ii)**, we use the PC algorithm to identify C from the inferred Z . The results in Table 3
661 demonstrate that the PC-adapted CiVAE is still significantly more robust to latent post-treatment bias
662 compared to CEVAE and TEDVAE, which empirically verify our claim that PC-adapted CiVAE can
663 address the case where inter-interactions exist among latent confounders and post-treatment variables.

664 **NeurIPS Paper Checklist**

665 **1. Claims**

666 Question: Do the main claims made in the abstract and introduction accurately reflect the
667 paper’s contributions and scope?

668 Answer: [\[Yes\]](#)

669 Justification: The contribution of this paper can be summarized as: We study a critical but
670 easily overlooked problem in causal effect estimation: latent post-treatment bias, and we
671 propose a novel framework, i.e., CiVAE, to address the bias. The details are in Section 4.

672 **2. Limitations**

673 Question: Does the paper discuss the limitations of the work performed by the authors?

674 Answer: [\[Yes\]](#)

675 Justification: We have discussed the potential issue of the vanilla when interactions among
676 the latent variables exists. However, in Section D we have addressed the issue by extending
677 our framework.

678 **3. Theory Assumptions and Proofs**

679 Question: For each theoretical result, does the paper provide the full set of assumptions and
680 a complete (and correct) proof?

681 Answer: [\[Yes\]](#)

682 Justification: We have introduced the three mild assumptions required for the identification
683 of causal effects under latent post-treatment bias. In addition, we have provided the proof
684 for all the theorems in the Appendix.

685 **4. Experimental Result Reproducibility**

686 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
687 perimental results of the paper to the extent that it affects the main claims and/or conclusions
688 of the paper (regardless of whether the code and data are provided or not)?

689 Answer: [\[Yes\]](#)

690 Justification: We have provided implementation details in Section 5.1. In addition, we have
691 provided a code demo in an anonymous URL.

692 **5. Open access to data and code**

693 Question: Does the paper provide open access to the data and code, with sufficient instruc-
694 tions to faithfully reproduce the main experimental results, as described in supplemental
695 material?

696 Answer: [\[Yes\]](#)

697 Justification: See Checklist 4.

698 **6. Experimental Setting/Details**

699 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
700 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
701 results?

702 Answer: [\[Yes\]](#)

703 Justification: See Checklist 4.

704 **7. Experiment Statistical Significance**

705 Question: Does the paper report error bars suitably and correctly defined or other appropriate
706 information about the statistical significance of the experiments?

707 Answer: [\[Yes\]](#)

708 Justification: We have reported the error of five independent run for both the proposed
709 CiVAE and all the baselines in the main paper.

710 **8. Experiments Compute Resources**

711 Question: For each experiment, does the paper provide sufficient information on the com-
712 puter resources (type of compute workers, memory, time of execution) needed to reproduce
713 the experiments?
714 Answer: [Yes]
715 Justification: See Checklist 4.

716 **9. Code Of Ethics**

717 Question: Does the research conducted in the paper conform, in every respect, with the
718 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>
719 Answer: [Yes]
720 Justification: We have carefully read the code of ethics and behaved strictly according to it.

721 **10. Broader Impacts**

722 Question: Does the paper discuss both potential positive societal impacts and negative
723 societal impacts of the work performed?
724 Answer: [Yes]
725 Justification: See Section A of the Appendix.

726 **11. Safeguards**

727 Question: Does the paper describe safeguards that have been put in place for responsible
728 release of data or models that have a high risk for misuse (e.g., pretrained language models,
729 image generators, or scraped datasets)?
730 Answer: [NA]
731 Justification: Our model does not have a high risk for misuse.

732 **12. Licenses for existing assets**

733 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
734 the paper, properly credited and are the license and terms of use explicitly mentioned and
735 properly respected?
736 Answer: [Yes]
737 Justification: We have cited the papers of our baselines and honor their license of code.

738 **13. New Assets**

739 Question: Are new assets introduced in the paper well documented and is the documentation
740 provided alongside the assets?
741 Answer: [Yes]
742 Justification: We provide the Readme file along side the codes.

743 **14. Crowdsourcing and Research with Human Subjects**

744 Question: For crowdsourcing experiments and research with human subjects, does the paper
745 include the full text of instructions given to participants and screenshots, if applicable, as
746 well as details about compensation (if any)?
747 Answer: [NA]
748 Justification: No human subjects are involved in our experiments.

749 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**
750 **Subjects**

751 Question: Does the paper describe potential risks incurred by study participants, whether
752 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
753 approvals (or an equivalent approval/review based on the requirements of your country or
754 institution) were obtained?
755 Answer: [NA]
756 Justification: No human subjects are involved in our experiments.