

## APPENDIX for “A Unified Causal View of Instruction Tuning”

### OVERVIEW:

- Appendix A contains the detailed data generating process, detailed proofs for all theoretical results in the main paper, as well as the proposed lemmas.
- Appendix B contains the prompt engineering to indicate task information.
- Appendix C contains the details of causal factor selection, including the encoding of the task representations and the mapping into latent mask vectors.
- Appendix D contains the details of causal factor constraint, including the key idea of the UIC Loss and the implementation of matrix  $A$  from the Theorem 2.
- Appendix E contains the details of tasks and datasets, including task selection and sampling strategy.
- Appendix F contains additional details of training and inference process.
- Appendix G contains additional experimental results under few-shot learning.

### A LEMMAS AND PROOFS

This section is structured as follows. We first provide some notations employed in this paper. In Appendix A.1, we provide a more detailed description for the data generating process in the main paper. Appendix A.2 presents the complete proof of Theorem 1. In Appendix A.3, we propose and prove useful lemmas that will be utilized in the proof of Theorem 2. Finally, Appendix A.4 offers the full proof of Theorem 2.

**Notations.** In this section, we adhere to a uniform notation scheme as in the main paper. Random variables are denoted by uppercase letters, while specific values are represented by lowercase letters, unless specified otherwise. For instance,  $\mathbf{X}$  is a random variable and  $\mathbf{x}$  is a particular value. Vector values are indicated by bold typeface (e.g.,  $\mathbf{x}$ ), while scalar values are represented using regular typeface (e.g.,  $x$ ). Additionally, calligraphic-style letters are used to denote representation spaces. For example,  $\mathcal{X}$  represents a representation space where  $\mathbf{x}$  belongs, with  $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^{\dim(\mathbf{x})}$ .

#### A.1 DATA GENERATING PROCESS

Before presenting the lemmas and proofs for identifiability, it is crucial to provide a comprehensive explanation of the data generating process. Understanding the data generating process is pivotal in the study of causality, as it unveils the causal mechanisms (denoted as assignment functions in Section 3.1) through which observed variables are produced by latent factors. In this regard, we employ the structural causal model (SCM), a widely utilized framework, to describe the data generating process. Formally, let  $\mathbf{x}_t \in \mathbb{R}^{\dim(\mathbf{x}_t)}$ ,  $\mathbf{y}_t \in \mathbb{R}^{\dim(\mathbf{y}_t)}$ ,  $\mathbf{l}_i \in \mathbb{R}^{\dim(\mathbf{l}_i)}$ . The parent set of  $\mathbf{X}_t$  denoted as  $Pa(\mathbf{X}_t)$  and the parent set of  $\mathbf{Y}_t$  denoted as  $Pa(\mathbf{Y}_t)$ . As explained in Section 3.1, the source context  $\mathbf{X}_t$  carries all the information of  $\mathbf{L}$ , hence  $Pa(\mathbf{X}_t) = \{\mathbf{L}_1, \mathbf{L}_2, \mathbf{L}_3, \dots, \mathbf{L}_n\}$ . In order to simplify the expression of exponential family distribution, we define  $\Theta_{\mathbf{x}_t} \triangleq \{f_{\mathbf{x}_t}, \Phi_{\mathbf{x}_t}\}$ , where  $f_{\mathbf{x}_t}$  denotes the invertible generating function,  $\Phi_{\mathbf{x}_t}$  represents the set of sufficient statistics  $\mathbf{T}$  and its coefficient  $\lambda$ .

The joint probability density of source context  $\mathbf{X}_t$  and latent factors  $\mathbf{L}_i$  can be written as:

$$p_{\Theta_{\mathbf{x}_t}}(\mathbf{x}_t, Pa(\mathbf{x}_t)|\mathbf{d}) = p_{\Theta_{\mathbf{x}_t}}(\mathbf{x}_t, Pa(\mathbf{x}_t)|\mathbf{d}) \quad (9)$$

$$= p_{f_{\mathbf{x}_t}}(\mathbf{x}_t|Pa(\mathbf{x}_t)) \cdot p_{\Phi_{\mathbf{x}_t}}(Pa(\mathbf{x}_t)|\mathbf{d}). \quad (10)$$

According to the additive noise model (ANM) assumption (Equation 1), the data generating process of  $\mathbf{x}_t$  can be written as:

$$\mathbf{x}_t = f_{\mathbf{x}_t}(Pa(\mathbf{x}_t)) + \varepsilon_{\mathbf{x}_t}, \quad \varepsilon_{\mathbf{x}_t} \sim p_{\varepsilon}(\varepsilon). \quad (11)$$

Using Equation 11, we can rewrite Equation 9 as:

$$p_{\Theta_{\mathbf{x}_t}}(\mathbf{x}_t, Pa(\mathbf{x}_t)|\mathbf{d}) = p_{f_{\mathbf{x}_t}}(\mathbf{x}_t|Pa(\mathbf{x}_t)) \cdot p_{\Phi_{\mathbf{x}_t}}(Pa(\mathbf{x}_t)|\mathbf{d}) \quad (12)$$

$$\Rightarrow p_{\Theta_{\mathbf{x}_t}}(\mathbf{x}_t, Pa(\mathbf{x}_t)|\mathbf{d}) = p_{\varepsilon_{\mathbf{x}_t}}(\mathbf{x}_t - f_{\mathbf{x}_t}(Pa(\mathbf{x}_t))) \cdot p_{\Phi_{\mathbf{x}_t}}(Pa(\mathbf{x}_t)|\mathbf{d}). \quad (13)$$

Considering that exponential family has universal approximation capability for probability density function, we assume the conditional probability density function  $p_{\Phi_{\mathbf{x}_t}}(Pa(\mathbf{x}_t)|\mathbf{d})$  is given by:

$$p_{\Phi_{\mathbf{x}_t}}(Pa(\mathbf{x}_t)|\mathbf{d}) = \prod_{i=1}^n p_{\mathbf{T}_i, \lambda_i}(\mathbf{l}_i|\mathbf{d}) \quad (14)$$

$$\Rightarrow p_{\Phi_{\mathbf{x}_t}}(Pa(\mathbf{x}_t)|\mathbf{d}) = \prod_{i=1}^n \prod_{j=1}^{\dim(\mathbf{l}_i)} p_{\mathbf{T}_i, \lambda_i}(l_{i,j}|\mathbf{d}) \quad (15)$$

$$\Rightarrow p_{\Phi_{\mathbf{x}_t}}(Pa(\mathbf{x}_t)|\mathbf{d}) = \prod_{i=1}^n \prod_{j=1}^{\dim(\mathbf{l}_i)} \frac{Q_{i,j}(l_{i,j})}{Z_{i,j}(\mathbf{d})} \exp \left[ \sum_{k=1}^{\dim(\mathbf{T}_{i,j})} T_{i,jk}(l_{i,j}) \lambda_{i,jk}(\mathbf{d}) \right]. \quad (16)$$

Notice that we employ a slightly different notation,  $p_{\mathbf{T}_i, \lambda_i}(\mathbf{l}_i|\mathbf{d})$ , instead of  $p_{\mathbf{L}_i}(\mathbf{l}_i|\mathbf{d})$ , to denote the conditional probability density of the latent factor  $\mathbf{l}_i$ , which is aimed at emphasizing that the latent factors are represented using exponential family distributions.

Equation 16 is called exponential family distribution, where  $Q_{i,j}$  is the base measure,  $Z_{i,j}$  is the partition function, i.e. normalization function,  $T_{i,jk}$  is one of the sufficient statistics and  $\lambda_{i,jk}$  is the corresponding coefficient. We can also rewrite  $T_{i,jk}$  and  $\lambda_{i,jk}$  in vector form:

$$\mathbf{T}_{i,j}(l_{i,j}) = [T_{i,j1}(l_{i,j}), T_{i,j2}(l_{i,j}), \dots, T_{i,jk}(l_{i,j})]^T. \quad (17)$$

$$\boldsymbol{\lambda}_{i,j}(\mathbf{d}) = [\lambda_{i,j1}(\mathbf{d}), \lambda_{i,j2}(\mathbf{d}), \dots, \lambda_{i,jk}(\mathbf{d})]^T. \quad (18)$$

Substituting it in Equation 16:

$$p_{\Phi_{\mathbf{x}_t}}(Pa(\mathbf{x}_t)|\mathbf{d}) = \prod_{i=1}^n \prod_{j=1}^{\dim(\mathbf{l}_i)} \frac{Q_{i,j}(l_{i,j})}{Z_{i,j}(\mathbf{d})} \exp [\boldsymbol{\lambda}_{i,j}(\mathbf{d})^\top \mathbf{T}_{i,j}(l_{i,j})]. \quad (19)$$

In this work, we adopt the following mild assumptions for the data generating processes, which are commonly used in other works (Khemakhem et al., 2020; Sun et al., 2021; Lu et al., 2021):

**Assumption 1 (Bijective).** *The generating functions  $f_{\mathbf{X}_t}$ ,  $f_{\mathbf{Y}_t}$  are bijective.*

**Assumption 2 (Denoising).** *Characterisitic functions of  $\varepsilon_{\mathbf{X}_t}$ ,  $\varepsilon_{\mathbf{Y}_t}$  are nonzero almost everywhere.*

**Assumption 3 (Transformation).** *The sufficient statistics  $\mathbf{T}$  are linear independent on every nonzero measure subset of  $L$  and are differentiable almost everywhere.*

**Assumption 4 (Variety).** *The number of different datasets, with differing inherent properties  $D$ , be  $n_D \geq n_0 = \max(\dim(\mathbf{l}_i) \times \dim(\mathbf{T}_{i,j})) + 1$ , and the following matrix has full column rank:*

$$\mathbf{H}_t = [\boldsymbol{\lambda}(\mathbf{d}_1) - \boldsymbol{\lambda}(\mathbf{d}_0), \boldsymbol{\lambda}(\mathbf{d}_2) - \boldsymbol{\lambda}(\mathbf{d}_0), \dots, \boldsymbol{\lambda}(\mathbf{d}_{n_0}) - \boldsymbol{\lambda}(\mathbf{d}_0)]. \quad (4)$$

Note that Assumption 1 is commonly used in identifiability works. Assumption 2 is generally satisfied for most continuous random variables, including Gaussian, exponential, and beta distributions. By applying Fourier transformation, this assumption helps eliminate the effect of noise in Equation 1. Assumption 3 is satisfied for all distributions belonging to the strongly exponential distribution family. Assumption 4 stipulates that the training datasets should contain a sufficient number of different datasets, and the full column rank of  $\mathbf{H}_t$  indicates that datasets should be diverse enough.

## A.2 PROOF OF THEOREM 1

**Theorem 1.** *Considering the data generating process described in Section 3.1, where  $\mathbf{X}_t$ ,  $\mathbf{Y}_{t, t \in \{t_1, t_2, \dots, t_m\}}$  are generated according to Equation 1, and  $\mathbf{L}_{i, i \in \{1, 2, \dots, n\}}$  has the distribution specified in Equation 2, as well as the fulfillment of Assumptions 1 - 4. We introduce a set of sets  $\mathcal{F}$  that describes the topology structure of a SCM and can be used to determine whether the SCM is identifiable.  $\mathcal{F}$  is generated by the following steps:*

1.  $\emptyset, Pa(\mathbf{X}_{t_1}), Pa(\mathbf{Y}_{t_1}), \dots, Pa(\mathbf{X}_{t_m}), Pa(\mathbf{Y}_{t_m}) \in \mathcal{F}$
2. Set  $\mathcal{A}, \mathcal{B} \in \mathcal{F} \Rightarrow \text{Set } \mathcal{A} - \mathcal{B}, \mathcal{B} - \mathcal{A} \in \mathcal{F}$ . Here  $\mathcal{A} - \mathcal{B} = \mathcal{A} \cap \bar{\mathcal{B}}$

The SCM is  $\sim_P$  identifiable if the set of sets  $\mathcal{F}$  includes all singleton sets  $\mathbf{L}_i$ , that is

$$\{\mathbf{L}_1\}, \{\mathbf{L}_2\}, \dots, \{\mathbf{L}_n\} \in \mathcal{F}.$$

*Proof.* The proof of the theorem can be roughly divided into two main steps. First, we transform the equations of probability density into an additive form. This step allows us to express the equations as a sum of individual components. Second, we apply the subtraction operator to the additive form equations, yielding equations with fewer latent factors. Consequently, each final equation contains only one of the latent factors.

**Step 1. Transforming** We begin our proof by stating that the learning marginal probability density on  $\mathbf{X}_t$  and  $\mathbf{Y}_t$  equals the true marginal probability density. For source context  $\mathbf{X}_t$ :

$$p_{\Theta_{x_t}}(\mathbf{x}_t) = p_{\bar{\Theta}_{x_t}}(\mathbf{x}_t) \quad (20)$$

$$\Rightarrow p_{f_{x_t}, \Phi_{x_t}}(\mathbf{x}_t | \mathbf{d}) = p_{\tilde{f}_{x_t}, \tilde{\Phi}_{x_t}}(\mathbf{x}_t | \mathbf{d}) \quad (21)$$

$$\begin{aligned} &\Rightarrow \int p_{f_{x_t}}(\mathbf{x}_t | Pa(\mathbf{x}_t)) p_{\Phi_{x_t}}(Pa(\mathbf{x}_t) | \mathbf{d}) \prod_{i=1}^n dl_i \\ &= \int p_{\tilde{f}_{x_t}}(\mathbf{x}_t | Pa(\mathbf{x}_t)) p_{\tilde{\Phi}_{x_t}}(Pa(\mathbf{x}_t) | \mathbf{d}) \prod_{i=1}^n dl_i \end{aligned} \quad (22)$$

$$\begin{aligned} &\Rightarrow \int p_{\varepsilon_{x_t}}(\mathbf{x}_t - f_{x_t}(Pa(\mathbf{x}_t))) p_{\Phi_{x_t}}(Pa(\mathbf{x}_t) | \mathbf{d}) \prod_{i=1}^n dl_i \\ &= \int p_{\varepsilon_{x_t}}(\mathbf{x}_t - \tilde{f}_{x_t}(Pa(\mathbf{x}_t))) p_{\tilde{\Phi}_{x_t}}(Pa(\mathbf{x}_t) | \mathbf{d}) \prod_{i=1}^n dl_i \end{aligned} \quad (23)$$

$$\begin{aligned} &\Rightarrow \int p_{\varepsilon_{x_t}}(\mathbf{x}_t - \bar{\mathbf{x}}_t) p_{\Phi_{x_t}}(f_{x_t}^{-1}(\bar{\mathbf{x}}_t) | \mathbf{d}) \left| \det(J_{f_{x_t}^{-1}}(\bar{\mathbf{x}}_t)) \right| d\bar{\mathbf{x}}_t \\ &= \int p_{\varepsilon_{x_t}}(\mathbf{x}_t - \bar{\mathbf{x}}_t) p_{\tilde{\Phi}_{x_t}}(\tilde{f}_{x_t}^{-1}(\bar{\mathbf{x}}_t) | \mathbf{d}) \left| \det(J_{\tilde{f}_{x_t}^{-1}}(\bar{\mathbf{x}}_t)) \right| d\bar{\mathbf{x}}_t \end{aligned} \quad (24)$$

$$\Rightarrow \int p_{\varepsilon}(\mathbf{x}_t - \bar{\mathbf{x}}_t) p_{\Phi_{x_t}, f_{x_t}, t}(\bar{\mathbf{x}}_t) d\bar{\mathbf{x}}_t = \int p_{\varepsilon}(\mathbf{x}_t - \bar{\mathbf{x}}_t) p_{\tilde{\Phi}_{x_t}, \tilde{f}_{x_t}, t}(\bar{\mathbf{x}}_t) d\bar{\mathbf{x}}_t \quad (25)$$

$$\Rightarrow (p_{\varepsilon_{x_t}} * p_{\Phi_{x_t}, f_{x_t}, t})(\mathbf{x}_t) = (p_{\varepsilon_{x_t}} * p_{\tilde{\Phi}_{x_t}, \tilde{f}_{x_t}, t})(\mathbf{x}_t) \quad (26)$$

$$\Rightarrow F[p_{\varepsilon_{x_t}}](\omega) F[p_{\Phi_{x_t}, f_{x_t}, t}](\omega) = F[p_{\varepsilon_{x_t}}](\omega) F[p_{\tilde{\Phi}_{x_t}, \tilde{f}_{x_t}, t}](\omega) \quad (27)$$

$$\Rightarrow F[p_{\Phi_{x_t}, f_{x_t}, t}](\omega) = F[p_{\tilde{\Phi}_{x_t}, \tilde{f}_{x_t}, t}](\omega) \quad (28)$$

$$\Rightarrow p_{\Phi_{x_t}, f_{x_t}, t}(\mathbf{x}_t) = p_{\tilde{\Phi}_{x_t}, \tilde{f}_{x_t}, t}(\mathbf{x}_t). \quad (29)$$

From Equation 21 to Equation 22, we introduce variables  $Pa(\mathbf{x}_t)$  into the formula and integrate them. This step is a commonly used technique to incorporate target variables in probability density equations. In Equation 24, the symbol  $J$  represents the Jacobian matrix, while  $|\det|$  denotes the generalized determinant of the matrix,  $\det|A| = \sqrt{\det(A^T A)}$ . In Equation 25, we introduce  $p_{\Phi_{x_t}, f_{x_t}, t}(\bar{\mathbf{x}}_t) = p_{\tilde{\Phi}_{x_t}}(\tilde{f}_{x_t}^{-1}(\bar{\mathbf{x}}_t) | \mathbf{d}) \left| \det(J_{\tilde{f}_{x_t}^{-1}}(\bar{\mathbf{x}}_t)) \right|$  for convenience. It is obviously that the Equation 25 is in the form of convolution. In Equation 26,  $F$  means Fourier transformation which is a useful tool to simplify convolution. From Equation 26 to Equation 28, we make an assumption that the characteristic function of noise  $F[p_{\varepsilon}]$  is non-zero almost everywhere, hence this term can be eliminated. Finally, we acquire the denoised result. Then taking the logarithm on the both sides of Equation 29 and

substituting the  $p_{\Phi_{\mathbf{x}_t}}$  with the exponential family distribution, we have

$$\begin{aligned}
& \log \left| \det(J_{f_{\mathbf{x}_t}^{-1}}(\mathbf{x}_t)) \right| \\
& + \sum_{i=1}^n \sum_{j=1}^{\dim(\mathbf{l}_i)} \left( Q_{i,j} \left( [f_{\mathbf{x}_t}^{-1}(\mathbf{x}_t)]_{i,j} \right) - Z_{i,j}(\mathbf{d}) + \sum_{k=1}^{\dim(\mathbf{T}_{i,j})} T_{i,jk} \left( [f_{\mathbf{x}_t}^{-1}(\mathbf{x}_t)]_{i,j} \right) \lambda_{i,jk}(\mathbf{d}) \right) \\
& = \log \left| \det(J_{\tilde{f}_{\mathbf{x}_t}^{-1}}(\mathbf{x}_t)) \right| \\
& + \sum_{i=1}^n \sum_{j=1}^{\dim(\mathbf{l}_i)} \left( \tilde{Q}_{i,j} \left( [\tilde{f}_{\mathbf{x}_t}^{-1}(\mathbf{x}_t)]_{i,j} \right) - \tilde{Z}_{i,j}(\mathbf{d}) + \sum_{k=1}^{\dim(\mathbf{T}_{i,j})} \tilde{T}_{i,jk} \left( [\tilde{f}_{\mathbf{x}_t}^{-1}(\mathbf{x}_t)]_{i,j} \right) \tilde{\lambda}_{i,jk}(\mathbf{d}) \right). \tag{30}
\end{aligned}$$

Notice that we have sufficient different tasks or datasets  $\mathbf{t}$ , that is, there exists  $\dim(\mathbf{l}_i) \times \dim(\mathbf{T}_{i,j}) + 1$  different  $t$ . Plugging these different  $\mathbf{t}$  in Equation 30 resulting to  $\dim(\mathbf{l}_i) \times \dim(\mathbf{T}_{i,j}) + 1$  equations. By subtracting the first equation from the second equation up to the last equation, we obtain a set of equations indexed by  $l = 1, 2, \dots, \dim(\mathbf{l}_i) \times \dim(\mathbf{T}_{i,j})$ :

$$\begin{aligned}
& \sum_i^n \left[ \langle \mathbf{T}_i \left( [f_{\mathbf{x}_t}^{-1}(\mathbf{x}_t)]_i \right), \bar{\lambda}_i(\mathbf{d}_l) \rangle + \sum_j \log \frac{Z_{i,j}(\mathbf{d}_0)}{Z_{i,j}(\mathbf{d}_l)} \right] \\
& = \sum_i^n \left[ \langle \tilde{\mathbf{T}}_i \left( [\tilde{f}_{\mathbf{x}_t}^{-1}(\mathbf{x}_t)]_i \right), \bar{\lambda}_i(\mathbf{d}_l) \rangle + \sum_j \log \frac{\tilde{Z}_{i,j}(\mathbf{d}_0)}{\tilde{Z}_{i,j}(\mathbf{d}_l)} \right]. \tag{31}
\end{aligned}$$

In Equation 31, we define  $\bar{\lambda}_i(\mathbf{d}_l) = \lambda_i(\mathbf{d}_l) - \lambda_i(\mathbf{d}_0)$ . In order to simplified Equation 31 further, we define  $\mathbf{w}_{l,i} = \sum_j \frac{\tilde{Z}_{i,j}(\mathbf{d}_0) Z_{i,j}(\mathbf{d}_l)}{\tilde{Z}_{i,i}(\mathbf{d}_l) Z_{i,j}(\mathbf{d}_0)}$ . Then we rewrite these equations in matrix form:

$$\sum_i^n \mathbf{H}_d^{i,\top} \mathbf{T}_i \left( [f_{\mathbf{x}_t}^{-1}(\mathbf{x}_t)]_i \right) = \sum_i^n \tilde{\mathbf{H}}_t^{i,\top} \tilde{\mathbf{T}}_i \left( [\tilde{f}_{\mathbf{x}_t}^{-1}(\mathbf{x}_t)] \right) + \mathbf{w}_{l,i}, \tag{32}$$

where  $\mathbf{H}_d^i = [\lambda_i(\mathbf{d}_1) - \lambda_i(\mathbf{d}_0), \lambda_i(\mathbf{d}_2) - \lambda_i(\mathbf{d}_0), \dots, \lambda_i(\mathbf{d}_{n_0}) - \lambda_i(\mathbf{d}_0)]$ ,  $n_0 = \dim(\mathbf{l}_i) \times \dim(\mathbf{T}_{i,j})$ .

**Step 2. Separation** Similar to  $\mathbf{x}_t$ , we can express the transformed equations for  $\mathbf{y}_t$  as well. Notice that the parent sets of  $\mathbf{x}_t$  encompass all latent factors  $\mathbf{l}_i$ , while the parent sets of  $\mathbf{y}_t$  usually encompass a subset of latent factors  $\mathbf{l}_i$ . We use the notation  $idx(Pa(\mathbf{Y}_t))$  to represent the indices of the latent factors comprising the set  $Pa(\mathbf{Y}_t)$ . We obtain  $m$  transformed equations for each  $\mathbf{Y}_{t_s}$ ,  $s = 1, 2, 3, \dots, m$ :

$$\sum_{i \in idx(Pa(\mathbf{Y}_{t_s}))} \mathbf{H}_t^{i,\top} \mathbf{T}_i \left( [f_{\mathbf{y}_{t_s}}^{-1}(\mathbf{y}_{t_s})]_i \right) = \sum_{i \in idx(Pa(\mathbf{Y}_{t_s}))} \tilde{\mathbf{H}}_t^{i,\top} \tilde{\mathbf{T}}_i \left( [\tilde{f}_{\mathbf{y}_{t_s}}^{-1}(\mathbf{y}_{t_s})] \right) + \mathbf{w}_{l,i}. \tag{33}$$

Furthermore, it is crucial to note that the latent factors  $\mathbf{l}_i$  are shared by  $\mathbf{Y}_t$ . Based on this property, we can express the transformed equations for the pair of target variables  $(\mathbf{y}_{t_s}, \mathbf{y}_{t_{s'}})$  as follows:

$$\sum_{\substack{i \in idx(Pa(\mathbf{Y}_{t_s}) \cup \\ Pa(\mathbf{Y}_{t_{s'}}))}} \mathbf{H}_t^{i,\top} \mathbf{T}_i \left( [f_{\mathbf{y}_{t_s}}^{-1}(\mathbf{y}_{t_s}, \mathbf{y}_{t_{s'}})]_i \right) = \sum_{\substack{i \in idx(Pa(\mathbf{Y}_{t_s}) \cup \\ Pa(\mathbf{Y}_{t_{s'}}))}} \tilde{\mathbf{H}}_t^{i,\top} \tilde{\mathbf{T}}_i \left( [\tilde{f}_{\mathbf{y}_{t_s}}^{-1}(\mathbf{y}_{t_s}, \mathbf{y}_{t_{s'}})] \right) + \mathbf{w}_{l,i}. \tag{34}$$

Notice that the subtraction of the two sets satisfies the following equation:

$$\mathcal{A} - \mathcal{B} \triangleq \mathcal{A} \cap \bar{\mathcal{B}} = (\mathcal{A} \cap \bar{\mathcal{B}}) \cup (\mathcal{B} \cap \bar{\mathcal{B}}) = (\mathcal{A} \cup \mathcal{B}) \cap \bar{\mathcal{B}} = (\mathcal{A} \cup \mathcal{B}) - \mathcal{B}. \tag{35}$$

Due to the inclusion property  $\mathcal{B} \subset \mathcal{A} \cup \mathcal{B}$ , the expression  $(\mathcal{A} \cup \mathcal{B}) - \mathcal{B}$  represents the removal of identical elements from the set  $\mathcal{A} \cup \mathcal{B}$  that are also present in  $\mathcal{B}$ . It is noteworthy that this type of set

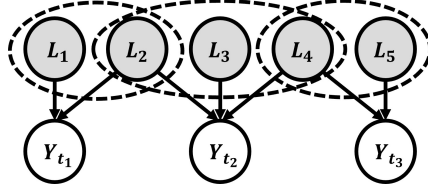


Figure 3: Identifiable latent factors.

subtraction demonstrates a striking similarity to algebraic subtraction. In parallel with the expansion of the set of sets  $\mathcal{F}$  through set subtraction, we can utilize algebraic subtraction on Equations 33 and Equations 34 to derive new equations that involve fewer latent factors. Given the condition that  $\mathcal{D}$  encompasses all singleton sets, it follows that all the latent factors can ultimately be isolated in their respective equations, as shown below:

$$\mathbf{H}_t^{i,\top} \mathbf{T}_i \left( \left[ f_{\mathbf{y}_{t_i}}^{-1}(\mathbf{y}_{t_s}) \right]_i \right) = \tilde{\mathbf{H}}_t^{i,\top} \tilde{\mathbf{T}}_i \left( \left[ \tilde{f}_{\mathbf{y}_{t_s}}^{-1}(\mathbf{y}_{t_s}) \right] \right) + \mathbf{w}_{l,i},$$

$$i \in \{1, 2, \dots, n\}, \quad \mathbf{t}_s \in \{\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_m\}. \quad (36)$$

Notice that the matrix  $\mathbf{H}_t^i$  has full rank, we multiply it's inverse matrix on both sides of Equation 36:

$$\mathbf{T}_i \left( \left[ f_{\mathbf{y}_{t_i}}^{-1}(\mathbf{y}_{t_s}) \right]_i \right) = \mathbf{M}_t^{i,\top} \tilde{\mathbf{T}}_i \left( \left[ \tilde{f}_{\mathbf{y}_{t_s}}^{-1}(\mathbf{y}_{t_s}) \right] \right) + \mathbf{v}_{l,i},$$

$$i \in \{1, 2, \dots, n\}, \quad \mathbf{t}_s \in \{\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_m\}, \quad (37)$$

where  $\mathbf{M}_t^i = (\mathbf{H}_t^{i,\top})^{-1} \tilde{\mathbf{H}}_t^{i,\top}$ ,  $\mathbf{v}_{l,i} = (\mathbf{H}_t^{i,\top})^{-1} \mathbf{w}_{l,i}$ .

Finally, we will prove that the matrix  $\mathbf{M}_t^i$  is a permutation matrix, demonstrating the  $\sim_P$  identifiability of the SCM. We adopt the method from Khemakhem et al. (2020) for this proof. Firstly, we consider the matrix  $\mathbf{T}$ . Under Assumption 4, the Jacobian of  $\mathbf{T}_i$  has a full column rank  $n$ , implying that the Jacobian of  $\mathbf{T}_i(f^{-1})$  is also of rank  $n$ . Consequently, the matrix  $\mathbf{M}_t^i$  is also of rank  $n$ . Secondly, we analyze two cases based on the dimension  $k$  of the sufficient statistics: (1)  $k = 1$ ; (2)  $k > 1$ . In the case of  $k = 1$ , the matrix  $\mathbf{T}_i$  becomes an  $n \times n$  square matrix. Since  $\mathbf{T}_i$  has a full rank, the matrix  $\mathbf{M}_t^i$  is also of full rank, indicating its invertibility. In the case of  $k > 1$ , we can directly apply Lemma 3 from Khemakhem et al. (2020) to prove the invertibility of  $\mathbf{M}_t^i$ . Lastly, assuming that both  $f$  and the sufficient statistics  $\mathbf{T}_i$  are twice differentiable, we apply Theorem 2 and Theorem 3 from Khemakhem et al. (2020) to demonstrate that  $\mathbf{M}_t^i$  is a permutation matrix.  $\square$

**Intuition.** To provide an intuitive understanding of Theorem 1, we present an identification process for Figure 3. Initially, we consider  $\mathbf{Y}_{t_1}$ , which is pointed by  $\mathbf{L}_1$  and  $\mathbf{L}_2$ . Solely relying on the information from  $\mathbf{Y}_{t_1}$  can not identify these latent factors. Next, we incorporate  $\mathbf{Y}_{t_2}$  into the analysis. By leveraging the information of  $\mathbf{Y}_{t_2}$ , we can identify  $\mathbf{L}_1$  and  $\mathbf{L}_2$ , for  $\mathbf{L}_1$  exclusively points to  $\mathbf{Y}_{t_1}$ , while  $\mathbf{L}_2$  points to both  $\mathbf{Y}_{t_1}$  and  $\mathbf{Y}_{t_2}$ . Subsequently, we include  $\mathbf{Y}_{t_3}$  in our analysis. Following the same procedure as before, the remaining three latent factors can be identified.

### A.3 LEMMAS

Before presenting the complete proof of Theorem 2, we first provide several useful lemmas.

**Lemma 1.** *Considering the data generating process described in Section 3.1. If there exist two distinct latent factors  $\mathbf{L}_i$  and  $\mathbf{L}_j$  such that their child sets  $Ch(\mathbf{L}_i)$  and  $Ch(\mathbf{L}_j)$  are identical, i.e.,  $Ch(\mathbf{L}_i) = Ch(\mathbf{L}_j)$ , then  $\mathbf{L}_i$  and  $\mathbf{L}_j$  can not be identified.*

*Proof.* We begin the proof with the equation of joint probability density:

$$p(\mathbf{X}_{\mathbf{t}_1}, \mathbf{Y}_{\mathbf{t}_1}, \dots, \mathbf{X}_{\mathbf{t}_m}, \mathbf{Y}_{\mathbf{t}_m} | \mathbf{d})$$

$$= p(\mathbf{X}_{\mathbf{t}_1}, \mathbf{Y}_{\mathbf{t}_1}, \dots, \mathbf{X}_{\mathbf{t}_m}, \mathbf{Y}_{\mathbf{t}_m} | \mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_n) \cdot p(\mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_n | \mathbf{d}) \quad (38)$$

$$= \prod_{t=\mathbf{t}_1}^{\mathbf{t}_m} p(\mathbf{X}_t | \mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_n) \cdot \prod_{t=\mathbf{t}_1}^{\mathbf{t}_m} p(\mathbf{Y}_t | Pa(\mathbf{Y}_t)) \cdot p(\mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_n | \mathbf{d}). \quad (39)$$

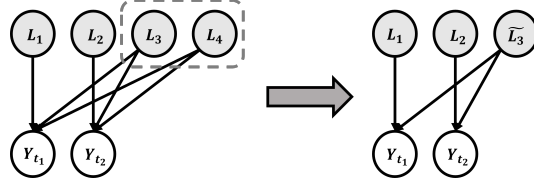


Figure 4: Unidentifiable latent factors.

We denoted  $Ch(\mathbf{L}_i) = Ch(\mathbf{L}_j) \triangleq Ch$ ,  $Ch = \{\mathbf{X}_{t_1}, \mathbf{X}_{t_2}, \dots, \mathbf{X}_{t_m}, \mathbf{Y}_{t'_1}, \dots, \mathbf{Y}_{t'_q}\}$ , in which  $\{\mathbf{Y}_{t'_1}, \dots, \mathbf{Y}_{t'_q}\} \subseteq \{\mathbf{Y}_{t_1}, \mathbf{Y}_{t_2}, \dots, \mathbf{Y}_{t_m}\}$ .

Back to the Equation 39,

$$\begin{aligned} & p(\mathbf{X}_{t_1}, \mathbf{Y}_{t_1}, \dots, \mathbf{X}_{t_m}, \mathbf{Y}_{t_m} | \mathbf{d}) \\ &= \prod_{t=t_1}^{t_m} p(\mathbf{X}_t | \mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_n) \cdot \prod_{t=t_1}^{t_m} p(\mathbf{Y}_t | Pa(\mathbf{Y}_t)) \cdot p(\mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_n | \mathbf{d}) \end{aligned} \quad (40)$$

$$\begin{aligned} &= \prod_{t=t_1}^{t_m} p(\mathbf{X}_t | (\mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_{-i}, \mathbf{L}_{-j}, \dots, \mathbf{L}_n), (\mathbf{L}_i, \mathbf{L}_j)) \\ &\cdot \prod_{t \in \{t'_1, \dots, t'_q\}} p(\mathbf{Y}_t | (Pa(\mathbf{Y}_t), \mathbf{L}_{-i}, \mathbf{L}_{-j}), (\mathbf{L}_i, \mathbf{L}_j)) \cdot \prod_{t \notin \{t'_1, \dots, t'_q\}} p(\mathbf{Y}_t | Pa(\mathbf{Y}_t)) \end{aligned} \quad (41)$$

$$\cdot p((\mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_{-i}, \mathbf{L}_{-j}, \dots, \mathbf{L}_n), (\mathbf{L}_i, \mathbf{L}_j) | \mathbf{d}). \quad (42)$$

Note that  $\mathbf{L}_i$  and  $\mathbf{L}_j$  always appear together in a term. Considering the following transformation:

$$(\mathbf{L}_i, \mathbf{L}_j) \rightarrow (\mathbf{L}'_i, \mathbf{L}'_j), \quad n = \min \left( \left\lfloor \frac{\dim(\mathbf{L}_i)}{2} \right\rfloor, \left\lfloor \frac{\dim(\mathbf{L}_j)}{2} \right\rfloor \right) \quad (43)$$

$$\mathbf{L}'_i = \begin{cases} \mathbf{L}'_{i[1:n]} = \mathbf{L}_{j[1:n]} \\ \mathbf{L}'_{i[n+1:\dim(\mathbf{L}_i)]} = \mathbf{L}_{i[n+1:\dim(\mathbf{L}_i)]} \end{cases}, \quad \mathbf{L}'_j = \begin{cases} \mathbf{L}'_{j[1:n]} = \mathbf{L}_{i[1:n]} \\ \mathbf{L}'_{j[n+1:\dim(\mathbf{L}_j)]} = \mathbf{L}_{j[n+1:\dim(\mathbf{L}_j)]} \end{cases} \quad (44)$$

The purpose of this transformation is to interchange the 1st to  $n$ th dimensions of  $\mathbf{L}_i$  and  $\mathbf{L}_j$ . As a result, the transformed variables  $\mathbf{L}'_i$  and  $\mathbf{L}'_j$  incorporate the information from both  $\mathbf{L}_i$  and  $\mathbf{L}_j$ . Note that both the original pair  $(\mathbf{L}_i, \mathbf{L}_j)$  and the transformed pair  $(\mathbf{L}'_i, \mathbf{L}'_j)$  satisfy Equation 39, indicating that it is impossible to uniquely recover the original pair  $(\mathbf{L}_i, \mathbf{L}_j)$  without information mixing. Consequently,  $\mathbf{L}_i$  and  $\mathbf{L}_j$  are not identifiable.  $\square$

**Intuition.** Figure 4 provides an intuitive understanding of Lemma 1. As depicted in Figure 4, when two latent factors  $\mathbf{L}_3$  and  $\mathbf{L}_4$  share the same child set  $\{\mathbf{Y}_{t_1}, \mathbf{Y}_{t_2}\}$ , it is equivalent to considering these two latent factors as a single variable.

**Lemma 2.** *Assuming the number of observed variables  $\mathbf{Y}$  is  $m$ , if the number of hidden variables  $\mathbf{Z}$  is greater than  $2^m - 1$ , then the causal graph is unidentifiable.*

*Proof.* Lemma 2 can be derived straightforwardly from Lemma 1. The number of different non-empty subsets of  $\{\mathbf{Y}_{t_1}, \mathbf{Y}_{t_2}, \dots, \mathbf{Y}_{t_m}\}$  is given by

$$\sum_i^m C_m^i = C_m^1 + C_m^2 + \dots + C_m^m = 2^m - 1. \quad (45)$$

$\square$

**Intuition.** Although the proof for Lemma 2 is technically straightforward, its meaning is quite interesting. Intuitively, Lemma 2 highlights the necessity of an adequate number of observed variables to identify latent factors. In this work, these observed variables correspond to distinct tasks or diverse datasets.

**Lemma 3.** *Assuming the number of observed variables  $\mathbf{Y}$  is  $m$ , for any observed variables  $\mathbf{Y}_{t_i}$ , its parent set satisfies the following:*

$$|Pa(\mathbf{Y}_{t_i})| \leq 2^{m-1}. \quad (46)$$

*In Equation 46, the notation  $|A|$  denotes the cardinality of a set  $A$ . For a finite set, the cardinality represents the number of elements it contains.*

*Proof.* We present a proof by contradiction. Let us assume that the given condition is violated, i.e.,  $|Pa(\mathbf{Y}_{t_i})| \geq 2^{m-1} + 1$ , which implies that there are at least  $2^{m-1} + 1$  latent factors  $\mathbf{L}$  pointing to  $\mathbf{Y}_{t_i}$ . Considering that all the child sets of these latent factors contain  $\mathbf{Y}_{t_i}$ , the only difference lies in the remaining  $m - 1$  latent factors. According to Lemma 2, the number of different child sets is limited to  $2^{m-1} - 1 + 1 = 2^{m-1}$  (including the empty set). However, the parent set  $Pa(\mathbf{Y}_{t_i})$  contains at least  $2^{m-1} + 1$  latent factors, indicating that there must exist two different latent factors with the same child set. This contradicts the initial assumption of identifiable latent factors. Consequently, we conclude that the condition  $|Pa(\mathbf{Y}_{t_i})| \leq 2^{m-1}$  holds.  $\square$

**Lemma 4.** *Considering a set of sets  $\mathcal{F}$  that describes the topology structure of a SCM.  $\mathcal{F}$  is generated by the following steps:*

1.  $\emptyset, Pa(\mathbf{X}_{t_1}), Pa(\mathbf{Y}_{t_1}), \dots, Pa(\mathbf{X}_{t_m}), Pa(\mathbf{Y}_{t_m}) \in \mathcal{F}$
2. Set  $\mathcal{A}, \mathcal{B} \in \mathcal{F} \Rightarrow$  Set  $\mathcal{A} - \mathcal{B}, \mathcal{B} - \mathcal{A} \in \mathcal{F}$ . Here  $\mathcal{A} - \mathcal{B} = \mathcal{A} \cap \bar{\mathcal{B}}$

*The set of sets  $\mathcal{F}$  includes all singleton sets  $\mathbf{L}_i$ , that is  $\{\mathbf{L}_1\}, \{\mathbf{L}_2\}, \dots, \{\mathbf{L}_n\} \in \mathcal{F}$ , **if and only if** ( $\Leftrightarrow$ ), For any two distinct latent factors  $\mathbf{L}_i$  and  $\mathbf{L}_j$  in the SCM, their child sets are not identical.*

*Proof.* We will first prove the direction " $\Rightarrow$ " (i.e., "only if"). We present a proof by contradiction. Let us assume that there exists two distinct latent factors  $\mathbf{L}_i$  and  $\mathbf{L}_j$  that have the same child sets, denoted as  $Ch = \{\mathbf{X}_{t_1}, \mathbf{X}_{t_2}, \dots, \mathbf{X}_{t_m}, \mathbf{Y}_{t'_1}, \dots, \mathbf{Y}_{t'_q}\}$ , where  $\{\mathbf{Y}_{t'_1}, \dots, \mathbf{Y}_{t'_q}\} \subseteq \{\mathbf{Y}_{t_1}, \mathbf{Y}_{t_2}, \dots, \mathbf{Y}_{t_m}\}$ . Let  $\bar{C}h = \{\mathbf{X}_{t_1}, \mathbf{Y}_{t_1}, \dots, \mathbf{X}_{t_m}, \mathbf{Y}_{t_m}\} - Ch = \{\mathbf{Y}_{t'_{q+1}}, \mathbf{Y}_{t'_{q+2}}, \dots, \mathbf{Y}_{t'_m}\}$ .

Notice that the original set  $\mathcal{F} = \{\emptyset, Pa(\mathbf{X}_{t_1}), Pa(\mathbf{Y}_{t_1}), \dots, Pa(\mathbf{X}_{t_m}), Pa(\mathbf{Y}_{t_m})\}$  can be divided into two distinct partition based on the sets  $Ch$  and  $\bar{C}h$ . The sets in one partition,  $\{\emptyset, Pa(\mathbf{Y}_{t'_{q+1}}), \dots, Pa(\mathbf{Y}_{t'_m})\} \subset \mathcal{F}$  do not includes either  $\mathbf{L}_i$  or  $\mathbf{L}_j$ , while the sets in the other partition,  $\{Pa(\mathbf{X}_{t_1}), Pa(\mathbf{X}_{t_2}), \dots, Pa(\mathbf{X}_{t_m}), Pa(\mathbf{Y}_{t'_1}), \dots, Pa(\mathbf{Y}_{t'_q})\} \subset \mathcal{F}$ , contains both  $\mathbf{L}_i$  and  $\mathbf{L}_j$ . Therefore, when performing the set subtraction, the result set can either contains both  $\mathbf{L}_i$  and  $\mathbf{L}_j$ , or it can contains neither  $\mathbf{L}_i$  nor  $\mathbf{L}_j$ , both of which still belong to one of the partitions. Hence, it is impossible to generate the singleton  $\{\mathbf{L}_i\}$  and  $\{\mathbf{L}_j\}$ , thus contradicting the assumption "the set of sets  $\mathcal{F}$  includes all singleton sets". Consequently, we conclude that the direction  $\Rightarrow$  holds.

Next, we will prove the direction " $\Leftarrow$ " (i.e., "if"). To begin, let us introduce the property of set subtraction. Consider two sets, denoted as  $\mathcal{A}$  and  $\mathcal{B}$ . Performing set subtraction on these two sets yields three distinct sets:  $\mathcal{A} - \mathcal{B}$ ,  $\mathcal{B} - \mathcal{A}$ , and  $\mathcal{A} - (\mathcal{A} - \mathcal{B})$ . Notably,  $\mathcal{B} - (\mathcal{B} - \mathcal{A})$  is equal to  $\mathcal{A} - (\mathcal{A} - \mathcal{B})$ , thus obviating the need to introduce this particular set. Furthermore, it is obvious that  $(\mathcal{A} - \mathcal{B}) \cup (\mathcal{B} - \mathcal{A}) \cup (\mathcal{A} - (\mathcal{A} - \mathcal{B})) = \mathcal{A} \cup \mathcal{B}$ . And the cardinality of three generated new sets are constrained by:  $\min(|\mathcal{A} - \mathcal{B}|, |\mathcal{B} - \mathcal{A}|, |\mathcal{A} - (\mathcal{A} - \mathcal{B})|) \leq \frac{1}{2} \max(|\mathcal{A}|, |\mathcal{B}|)$ .

Let us now consider the original set  $\mathcal{F} = \{\emptyset, Pa(\mathbf{X}_{t_1}), Pa(\mathbf{Y}_{t_1}), \dots, Pa(\mathbf{X}_{t_m}), Pa(\mathbf{Y}_{t_m})\}$ . Note that the parent sets of every  $\mathbf{X}_t$  contain all the latent factors  $\mathbf{L}$ , which can be represented as the universal set  $\mathcal{U}$ . Therefore, we can select one of the  $\mathbf{X}_t$ , denoted as  $\mathbf{X}$ , as it encompasses the entire set of latent factors. Next, we consider a total of  $m + 1$  observed variables, where  $m$  of them are denoted as  $\mathbf{Y}_t$ , and one of them is  $\mathbf{X}$ . According to Lemma 3, the cardinality of their parent sets is no more than  $2^m$ . Here we present a set generating process: Firstly, we have a set  $Pa(\mathbf{X})$ . Next, by introducing the set  $Pa(\mathbf{Y}_{t_1})$  and performing set subtraction, we obtain three new sets  $Pa(\mathbf{X}) - Pa(\mathbf{Y}_{t_1})$ ,  $Pa(\mathbf{Y}_{t_1}) - Pa(\mathbf{X})$  and  $Pa(\mathbf{X}) - (Pa(\mathbf{X}) - Pa(\mathbf{Y}_{t_1}))$ . Subsequently, we introduce the set  $Pa(\mathbf{Y}_{t_2})$  and perform set subtraction on each of these three sets, resulting in nine new sets. We repeat this process by introducing  $Pa(\mathbf{Y}_{t_2}), \dots, Pa(\mathbf{Y}_{t_m})$  and performing set subtraction. Finally, we obtain  $3^m$  generated sets denoted as  $\mathcal{S}$ . As mentioned earlier, the union set of these  $3^m$  generated sets is the universal set  $\mathcal{U}$ . Moreover, the cardinality of these sets is constrained

by the following condition:

$$|\mathcal{S}| \leq \frac{1}{2} \cdot \frac{1}{2} \cdots \frac{1}{2} \cdot |Pa(X)| \leq \left(\frac{1}{2}\right)^m \cdot 2^m = 1. \quad (47)$$

Equation 47 indicates that the cardinality of each generated set is no more than 1, implying that they are either empty sets or singletons. Combining this with the fact that the union set is the universal set, we can conclude that  $\{\mathbf{L}_1\}, \{\mathbf{L}_2\}, \dots, \{\mathbf{L}_n\} \in \mathcal{F}$ . Therefore, the direction  $\Leftarrow$  holds.  $\square$

#### A.4 PROOF OF THEOREM 2

**Theorem 2.** *Considering the data generating process described in Section 3.1, we employ a binary adjacency matrix denoted as  $A$  to represent the topology relations between  $\mathbf{L}_i$  and  $\mathbf{Y}_t$ . The matrix  $A$  comprises  $m$  rows and  $n$  columns, where  $m$  is the number of  $\mathbf{Y}_t$ , and  $n$  is the number of latent factors  $\mathbf{L}_i$ . Specifically, a value of 1 at position  $(i, j)$  indicates that  $\mathbf{L}_j$  has a direct effect on  $\mathbf{Y}_i$ , while a value of 0 indicates no direct effect. Latent factors in a SCM are identifiable if and only if the following equation holds. We refer to the equation as the **uniform identifiability condition (UIC)**.*

$$\mathbb{1} \left( \frac{1}{m} [A^\top A + (1 - A)^\top (1 - A)] - I_{n \times n} \right) = 0_{n \times n}. \quad (5)$$

In Equation 5,  $\mathbb{1}(\cdot)$  is an indicator function, which defined as  $[\mathbb{1}(A)]_{ij} = \begin{cases} 0, & 0 \leq a_{ij} < 1 \\ 1, & a_{ij} = 1 \end{cases}$

*Proof.* The proof consists of three steps, and an overview of the proof is presented in Figure 5.

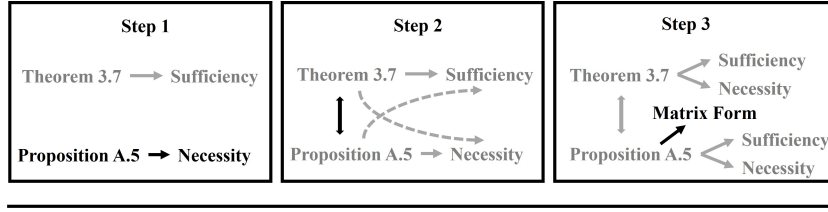


Figure 5: Overview of the proof. Each step focuses on the element marked in black. In Step 1, we demonstrate that the condition stated in Proposition 1 is a necessary condition for determining SCM identifiability. In Step 2, we establish the equivalence between the conditions in Proposition 1 and Theorem 1, thereby showing that both conditions are necessary and sufficient. Finally, in Step 3, we present the matrix form representation of the condition in Proposition 1.

**Step 1. Proving Necessity** We introduce a criterion to determine the identifiability of a given SCM. The criterion is that: For any two distinct latent factors  $\mathbf{L}_i$  and  $\mathbf{L}_j$  in the SCM, their child sets (i.e., sets containing  $\mathbf{X}_t$  and  $\mathbf{Y}_t$  pointed by  $\mathbf{L}$ ) are not identical. Then, according to Lemma 1, a negative answer to this criterion implies the non-identifiability of the SCM. Thus it can be seen as the contrapositive form of the necessary condition for identifiability. We can express the equivalent necessary condition in the form of a proposition:

**Proposition 1.** *If the SCM is identifiable, then for any two distinct latent factors  $\mathbf{L}_i$  and  $\mathbf{L}_j$  in the SCM, their child sets are not identical.*

**Step 2. Combining Necessity and Sufficiency** Notice that Theorem 1 provides a sufficient condition for the identifiability of SCM, while Proposition 1 presents a necessary condition for the identifiability of SCM. According to Lemma 4, these two conditions are exactly equivalent. Consequently, we conclude that both conditions are both necessary and sufficient for the identifiability of SCM. Based on the conclusion, we can strengthen Proposition 1 by incorporating the sufficiency aspect, as presented in Proposition 2.

**Proposition 2.** *A SCM is identifiable, if and only if for any two distinct latent factors  $\mathbf{L}_i$  and  $\mathbf{L}_j$  in the SCM, their child sets are not identical.*



**Step 3. Matrix Representation** In this step, we will represent the conditions using matrix notation. Notice that the condition described in Theorem 1 involves a generative process, which poses challenges when attempting to express it in matrix form. Therefore, we choose to employ the condition introduced in Proposition 2, i.e., for any two distinct latent factors  $\mathbf{L}_i$  and  $\mathbf{L}_j$  in the SCM, their child sets are not identical. This condition can be naturally expressed using a binary adjacency matrix denoted as  $A$ . The matrix  $A$  comprises  $m$  rows and  $n$  columns, where  $m$  is the number of  $\mathbf{Y}_t$ , and  $n$  is the number of latent factors  $\mathbf{L}_i$ . Specifically, a value of 1 at position  $(i, j)$  indicates that  $\mathbf{L}_j$  has a direct effect on  $\mathbf{Y}_i$ , while a value of 0 indicates no direct effect. The condition that the child sets are not identical is equivalent to stating that any two distinct columns in matrix  $A$  are not the same. Hence, we can express Proposition 2 in matrix form as Proposition 3.

**Proposition 3.** *Considering the binary adjacency matrix  $A$  described in Step 3, a SCM is identifiable, if and only if any two distinct columns in matrix  $A$  are not the same.*

Notice the following Equation 48 holds:

$$x_1 = \{0, 1\}, \quad x_2 = \{0, 1\}, \quad x_1x_2 + (1 - x_1)(1 - x_2) = \begin{cases} 1 & x_1 = x_2 \\ 0 & x_1 \neq x_2 \end{cases} \quad (48)$$

The formula  $x_1x_2 + (1 - x_1)(1 - x_2)$  can be regarded as a correlation function for  $x_1$  and  $x_2$ , and this correlation function can be straightforward generalized to a vector form:

$$C_{ij} \triangleq \text{Corr}(\mathbf{v}_i, \mathbf{v}_j) = \frac{1}{\dim(\mathbf{v})} [(\mathbf{v}_i^\top \mathbf{v}_j) + (1 - \mathbf{v}_i)^\top (1 - \mathbf{v}_j)], \quad (49)$$

where the term  $\frac{1}{\dim(\mathbf{v})}$  serves as a normalization factor.  $C_{ij} = 0$  if all of the elements in the same position of  $\mathbf{v}_i$  and  $\mathbf{v}_j$  are different.  $0 < C_{ij} < 1$  if some of the elements in the same position of  $\mathbf{v}_i$  and  $\mathbf{v}_j$  are the same.  $C_{ij} = 1$  if  $\mathbf{v}_i$  and  $\mathbf{v}_j$  are exactly the same.

Based on that, we can express the condition that "any two distinct columns in matrix  $A$  are not the same" using an equivalent matrix formula, as shown in Equation 5:

$$\mathbb{1} \left( \frac{1}{m} [A^\top A + (1 - A)^\top (1 - A)] - I_{n \times n} \right) = 0_{n \times n}.$$

Here the indicator function  $\mathbb{1}(\cdot)$  acts as a selector to identify which two columns are identical.  $\square$

## B PROMPT ENGINEERING

Table 5: Design of discrete prompt described in natural language. For classification tasks, we provide category options as part of prompt.

Task	Discrete Prompt
SUM	Summarize the document:
RC	Answer the question based on its following passage:
TC	Distinguish which topic the text is (options are [option]):
PD	Distinguish whether the two sentences have the same meaning (options are [option]):
SA	Distinguish which sentiment the review is (options are [option]):
LA	Distinguish whether the sentence is linguistically acceptable (options are [option]):
NLI	Distinguish whether the first sentence can infer its following sentence (options are [option]):

For both Vanilla-IT and SIT, we apply the same setting of prompt engineering as follow.

We adopt hybrid prompts  $p = \{p_d, p_c\}$  as instructions following (Xu et al., 2022; Chen et al., 2023), where discrete prompts  $p_d$  are natural words, while continuous prompts  $p_c$  are continuous embeddings. For the discrete prompts  $p_d$ , we manually design them as shown in Table 5. For the continuous prompts  $p_c$ , we utilize an individual prompt encoder to encode a sequence of trainable dense vectors. The prompt encoder is composed of two-layer bidirectional long-short term memory network (BiLSTM) (Graves & Graves, 2012) followed by a multilayer perceptron (MLP), i.e.,

$p_c = \text{MLP}(\text{BiLSTM}([p_1], [p_2], \dots, [p_{|p_c|}]))$ , where  $[p_j]_{j=1}^{|p_c|}$  represents placeholders to be replaced by trainable dense vectors, of length  $|p_c| = 6$  for each input sequence. Note that multiple source sequences are concatenated into one as input. In this work, there are at most two source sequences, and the prompted input is  $\langle p, x_1, x_2 \rangle = \{[s], p_d, p_c, x_1, [e], p_c, x_2, [e]\}$  for such tasks.

For the prompt encoder, the mid-hidden size and output size of the LSTM is 512 and 1024, respectively. Dropout with probability 0.1 is applied for LSTM. MLP is composed of two linear layers with a ReLU activation function in between. The hidden size and output size of the two-layer MLP is 1024.

## C DETAILS OF CAUSAL FACTOR SELECTION

In this section, we introduce the implementation details of the task representation  $\mathbf{h}_t$  and the latent mask vector  $\mathbf{m}_t$ .

**Task Representation.** We obtain task representation  $\mathbf{h}_t$  by encoding hybrid prompts  $p = \{p_d, p_c\}$  introduced in Appendix B. Specifically, for discrete prompts with variable length, we derive a single embedding  $\mathbf{e}_{p_d} \in \mathbb{R}^{d_h}$  through the utilization of average pooling, applied to the output embedding sequence generated from a word embedding layer. Also, for continuous prompts with the maximum length of 12 (the length twice as long as 6 for two source sequences), we linearly transform the output embedding sequence from the prompt encoder into another embedding  $\mathbf{e}_{p_c} \in \mathbb{R}^{d_h}$ . Then, we linearly combine them to achieve the task representation  $\mathbf{h}_t \in \mathbb{R}^n$ , i.e.,  $\mathbf{h}_t = \mathbf{W}_4 \mathbf{e}_{p_d} + \mathbf{W}_5 \mathbf{e}_{p_c} + \mathbf{b}_4$ .

**Latent Mask Vector.** We obtain the latent mask vector  $\mathbf{m}_t$  based on the task representation  $\mathbf{h}_t$ . Firstly,  $\mathbf{h}_t$  is normalized by a sigmoid activation function into  $\hat{\mathbf{h}}_t$ , a soft version of latent mask vector, i.e.,

$$\hat{\mathbf{h}}_t = \text{Sigmoid}(\mathbf{h}_t), \quad (50)$$

whose continuous value  $\hat{h}_{ti} \in (0, 1)$  in each dimension represents the selected probability of each latent factor. Then, we utilize bernoulli sampling to obtain the hard latent mask vector  $\mathbf{m}_t$  according to  $\hat{\mathbf{h}}_t$ , where the discrete value  $m_{ti} \in \{0, 1\}$  in each dimension is sampled from  $\{0, 1\}$  and only 1 represents "selected". To increase the stability of sampling results, we additionally multiply a scaling coefficient selected from (50, 200) for  $\mathbf{h}_t$  before the sigmoid activation.

## D DETAILS OF CAUSAL FACTOR CONSTRAINT

**UIC Loss Function.** Note that Equation 5 provides a necessary and sufficient condition for identifying latent factors. Using this equation, we can design a loss function to ensure identifiability in our model. However, a challenge arises with the indicator function in Equation 5, which is non-differentiable at  $a_{ij} = 1$ . This prevents direct application of gradient-based optimization methods. One solution is to replace the indicator function with an approximate, but differentiable function. In this work, we choose the power function  $x^\alpha$ , which becomes increasingly similar to the indicator function as  $\alpha$  approaches infinity (Practically,  $\alpha = 50$  is quite enough). Therefore, our loss function can be expressed as:

$$\mathcal{L}_{uic} = \frac{1}{m^\alpha} \sum_{i,j=1}^n \left[ \sum_{k=1}^m a_{ki} a_{kj} + (1 - a_{ki})(1 - a_{kj}) - m\delta_{ij} \right]^\alpha, \quad (6)$$

In Equation 6,  $\delta$  is the Kronecker delta function, which defined as  $\delta_{ij} = \begin{cases} 0, & i \neq j \\ 1, & i = j \end{cases}$

**Implementation of Matrix  $A$ .** To apply the UIC loss and task distinction loss, we construct a discrete task-latent matrix  $\mathbf{M}_{tl}$  to implement the binary adjacency matrix  $A$  described in Theorem 2, whose elements  $a$  are utilized in  $\mathcal{L}_{uic}$  (Equation 6) and  $\mathcal{L}_{dis}$  (Equation 7).

First, we construct a continuous version of this matrix. Specifically, we collect the soft latent mask vectors  $\hat{\mathbf{h}}_t \in \mathbb{R}^n$  (introduced in Appendix C) for  $m$  training mixture tasks, and stack  $m$  vectors into a continuous matrix  $\mathbf{M} \in \mathbb{R}^{m \times n}$ . Collected from training batches, this matrix changes dynamically.

Then, we discretize this continuous matrix. Since the bernoulli sampling does not meet the requirement of derivability, we apply gumbel-softmax rather than bernoulli sampling to realize discretization of  $\mathbf{M}_{t_l}$  for parameter optimization during training.

## E DETAILS OF TASKS AND DATASETS

In this section, we present the task selection as well as our sampling strategy.

**Task Selection.** We selected tasks based on established prior works like FLAN and T0, choosing widely-adopted datasets to validate our approach. This selection covers a diverse range of classical NLP tasks, including the General Language Understanding Evaluation (GLUE) benchmark, which is one of the most popular evaluation benchmarks used for multitask learning (Worsham & Kalita, 2020). Besides Natural Language Understanding (NLU) tasks, we also consider Natural Language Generation (NLG) tasks, e.g., Summarization and Reading Comprehension.

The training datasets consist of XSUM (Narayan et al., 2018), CNNDM (See et al., 2017), Duorc<sub>self</sub> (Saha et al., 2018), Duorc<sub>para</sub> (Saha et al., 2018), AG (Zhang et al., 2015), Trec (Li & Roth, 2002), PAWS (Zhang et al., 2019), IMDB (Maas et al., 2011) and Yelp (Zhang et al., 2015). The held-out datasets used are Gigaword (Napoles et al., 2012), Squad (Rajpurkar et al., 2016), DBPedia (Lehmann et al., 2015), MRPC (Dolan & Brockett, 2005), QQP (Wang et al., 2018), SST-2 (Socher et al., 2013), CoLA (Warstadt et al., 2019), MNLI<sub>m</sub> (Williams et al., 2018), MNLI<sub>mm</sub> (Williams et al., 2018), QNLI (Rajpurkar et al., 2018), RTE (Dagan et al., 2006), WNLI (Levesque et al., 2012). Details of all datasets are provided in Table 6a.

**Sampling Strategy.** To construct the training mixture dataset, we randomly sample and mix data from each dataset listed in Table 6a. Following the approach described in (Wei et al., 2021a; Raffel et al., 2020), we adopt the examples-proportional mixing scheme and limit the number of training examples per dataset to 15k. In order to increase the coverage of the sampled dataset with respect to the original dataset, we prioritize sampling data that has not been sampled before. Consequently, the sample size of the training mixture datasets in our work can be expressed as:

$$\text{size} = \min(\text{num}(\text{epochs}) \times 15\text{k}, \text{size}(\text{original dataset})), \quad (51)$$

where the number of training epochs is 10 in our works. The statistics of the final training mixture datasets and the held-out datasets are shown in Table 6.

Table 6: Data statistics.

(a) The training mixture datasets.				(b) The held-out datasets.				
Task	Dataset	Train (sampled)	Test	Task	Dataset	Split	Size	
SUM	XSUM	150,000	11,334	SUM	Gigaword	test	1,951	
	CNNDM	150,000	11,490	RC	Squad	dev	10,570	
RC	Duorc <sub>self</sub>	60,721	12,559	TC	DBPedia	test	70,000	
	Duorc <sub>para</sub>	69,524	15,857	PD	MRPC	dev	408	
TC	AG	120,000	7,600		QQP	dev	40,430	
	Trec	5,452	500	SA	SST-2	dev	872	
PD	PAWS	49,401	8,000	LA	CoLA	dev	1,043	
SA	IMDB	25,000	25,000	NLI	MNLI <sub>m</sub>	dev	9,815	
	Yelp	150,000	50,000			MNLI <sub>mm</sub>	dev	9,815
						QNLI	dev	5,463
						RTE	dev	277
						WNLI	dev	71

## F DETAILS OF TRAINING AND INFERENCE

In this section, we supplement more details about the training and inference process. For the tasks with one source sequence, we set the max length as 550, while for those with two source sequences, we set the max length as 350 for the first sentence, and 200 for the second sentence. For other hyper-parameters, we manually tune them based on the validation set or a subset of training set. Specifically, the batch size is selected from  $\{256, 512\}$ , the learning rate is selected from  $\{1e^{-5}, 3e^{-5}, 5e^{-5}\}$ . The total training steps contain 10 epochs, and we conduct evaluation for early stopping every epoch and every 500 steps. During inference, we apply beam search for text generation and set beam size as 6. Specifically, we use Huggingface Transformers library<sup>4</sup> for implementations<sup>5</sup>. All the reported results come from evaluation on models trained in the mixture datasets, which are subsets sampled from the full datasets.

## G FEW-SHOT LEARNING

In this section, we show all the experimental results under the few-shot setting in Table 7. The hyper-parameter setup is the same as the setup during training stage, except for the warm-up strategy absent in few-shot training. The last checkpoint are picked for prediction.

As shown in Table 7, SIT outperforms Vanilla-IT on 9 out of 12 datasets, demonstrating the better learning capability and generalization ability of SIT, which benefits from SCM capturing the underlying causal relationships. On the whole, the model performance improves more on the difficult tasks after few-shot learning, e.g. SUM task, while the performance on the simple tasks maybe decrease, e.g., RC task. We analyze the two cases in detail as follows. (i) On the datasets that have poor zero-shot performance, e.g., DBPedia and CoLA, both Vanilla-IT and SIT gain significantly under the few-shot setting as shown in Figure 2. The larger gain of SIT than Vanilla-IT indicates that structural instructions can adapt faster and better to a new target space with SCM as bridge between the task and target. (ii) On the datasets that have good zero-shot performance, e.g., SST-2, Vanilla-IT can only improve 0.87% in terms of accuracy by learning few samples, while SIT leads to a decrease in model performance. The possible reason is that with  $3e^{-5}$  as learning rate the same as training stage, the update rate of the model parameters is too fast, so that the prediction behavior is unstable or even worse for the tasks previously performed well. More suitable hyper-parameter setup needs to be determined by grid search.

Table 7: Few shot performance of all the held-out datasets, including OOD and cross-task situations.

Method	OOD Performance					Cross-task Performance						
	SUM	RC	TC	PD	SA	LA	NLI					
	Gigaword	Squad	DBPedia	MRPC	QQP	SST-2	CoLA	MNLI <sub>m</sub>	MNLI <sub>mm</sub>	QNLI	RTE	WNLI
Vanilla-IT	29.82	54.02	76.33	68.38	36.82	<b>93.23</b>	38.16	32.65	32.94	50.52	16.97	43.66
SIT	<b>30.05</b>	<b>75.99</b>	<b>93.16</b>	68.38	<b>74.52</b>	87.96	<b>69.03</b>	<b>35.39</b>	<b>35.21</b>	<b>50.54</b>	<b>47.29</b>	43.66

<sup>4</sup><https://github.com/huggingface/transformers>

<sup>5</sup>The code is available at <https://anonymous.4open.science/r/SIT-34DB/>