# I PREDICT THEREFORE I AM: IS NEXT TOKEN PREDICTION ENOUGH TO LEARN HUMAN-INTERPRETABLE CONCEPTS FROM DATA?

# **Anonymous authors**

Paper under double-blind review

#### **ABSTRACT**

Recent empirical evidence shows that LLM representations encode human-interpretable concepts. Nevertheless, the mechanisms by which these representations emerge remain largely unexplored. To shed further light on this, we introduce a novel generative model that generates tokens on the basis of such concepts formulated as latent discrete variables. Under mild conditions, even when the mapping from the latent space to the observed space is non-invertible, we establish rigorous identifiability result: the representations learned by LLMs through next-token prediction can be approximately modeled as the logarithm of the posterior probabilities of these latent discrete concepts given input context, up to an invertible linear transformation. This theoretical finding: 1) provides evidence that LLMs capture essential underlying generative factors, 2) offers a unified and principled perspective for understanding the linear representation hypothesis, and 3) motivates a theoretically grounded approach for evaluating sparse autoencoders. Empirically, we validate our theoretical results through evaluations on both simulation data and the Pythia, Llama, and DeepSeek model families.

### 1 Introduction

Large language models (LLMs) are trained on extensive datasets, primarily sourced from the Internet, enabling them to excel in a wide range of downstream tasks, such as language translation, text summarization, and question answering (Zhao et al., 2024; Bommasani et al., 2021; Olah et al., 2020). Despite this success, their internal representations remain largely opaque, a fact that has sparked a series of efforts toward theoretical understanding and practical interpretability. A promising direction arises from recent empirical evidence showing that LLM representations (often called activations in the natural language processing community) encode latent concepts, such as sentiment (Turner et al., 2023) or writing style (Lyu et al., 2023), that align with human-interpretable abstractions (Acerbi & Stubbersfield, 2023; Manning et al., 2020; Sajjad et al., 2022; Sharkey et al., 2025). Uncovering the underlying reasons for this phenomenon, i.e., developing a principled framework that formally links learned representations to latent concepts, holds promise for advancing our understanding of LLMs.

Some recent works attempt to formulate the problem of linking LLM representations to latent concepts through latent variable models (Park et al., 2023; 2024; Rajendran et al., 2024) in which human-interpretable concepts are formulated as latent variables. For example, Park et al. (2023) propose a latent variable model that focuses solely on binary latent concepts. The work of Park et al. (2024) extends this formalization to categorical concepts, but its focus on hierarchical relationships between concepts may not capture other possible structures in textual data. Rajendran et al. (2024) instead model both latent concepts and, in particular, observed text data, as continuous variables which may deviate from the discrete sense of language. See Appendix A for additional related work.

We investigate here how LLMs, despite relying solely on next-token prediction, can grasp humaninterpretable concepts, through a new latent variable model without the limitations above.

We begin by introducing a latent variable model, illustrated on the left of Figure 1. In this model, human-interpretable concepts are formulated as latent variables connected by arbitrary causal relationships, which in turn generate observed text data through an underlying generative process. Building on this model we develop a theoretical framework to assess whether next-token prediction enables

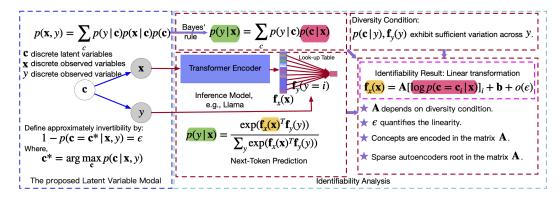


Figure 1: An overview of the main contributions of this work. On the left, we illustrate the proposed latent variable model that represents concepts as latent variables  $\mathbf{c}$ , which generates both the input  $\mathbf{x}$  and output y within a next-token prediction framework. Leveraging Bayes' rule and the diversity condition, we establish an identifiability result: LLM representations are related to a linear transformation of the logarithm of the posterior distribution of latent variables conditioned on input context, i.e.,  $\mathbf{f_x}(\mathbf{x}) = \mathbf{A}[\log p(\mathbf{c} = \mathbf{c}_i|\mathbf{x})]_i + \mathbf{b} + o(\epsilon)$ , where  $\mathbf{b}$  is a constant, and  $o(\epsilon)$  represents a term that grows asymptotically smaller than  $\epsilon$  as  $\epsilon \to 0$ .

LLMs trained on such text data to recover the underlying latent variables. Through identifiability analysis, we show that, under mild conditions, LLM representations approximately correspond to a *linear* transformation of the logarithm of the posterior distribution of latent variables conditioned on input context. We refer to this as the *linear property* of LLM representations.

Taking a step further, we demonstrate that this *linear property* provides both theoretical and practical insights for *linear* phenomena observed in empirical studies. On the theoretical side, it offers a unified perspective for understanding the various forms from the linear representation hypothesis, e.g., steering vector (Turner et al., 2023) and linear probing (Park et al., 2023; Marks & Tegmark, 2023). On the practical side, the linear property motivates a principled evaluation strategy for sparse autoencoders (SAEs), which aim to extract concepts in an *unsupervised* manner (Huben et al., 2024; Bricken et al., 2023). Specifically, *supervised* concept extraction methods, e.g., linear probing, can serve as an upper bound to evaluate how well SAE features recover the underlying concepts.

Contributions. To investigate the relationship between LLM representations and human-interpretable concepts, this paper makes the following contributions. We introduce a novel latent variable model in Sec. 2, in which human-interpretable concepts are formulated as latent variables. Building on this setup, we establish the existence of a *linear* transformation between LLM representations and human-interpretable concepts in Sec. 3. In Sec. 4, we show how the linear property above provides theoretical support for the empirical linear representation hypothesis in LLMs and establishes a unified perspective. In Sec. 5, we show how the *linear* property motivate a principled approach for evaluating SAEs. Additionally, we offer an early exploration of a promising method: structured SAEs, a variant that incorporates structured sparsity, to improve performance under this evaluation. In Sec. 6, we present experiments on simulation data, and real data across the Pythia (Biderman et al., 2023), Llama (Touvron et al., 2023a; Dubey et al., 2024), and DeepSeek-R1 (Guo et al., 2025) models families, with results consistent with our findings. We further present empirical results on the proposed evaluation approach, and the structured SAEs.

### 2 Setup: Latent Variable Model for Text Data Generation

We begin by a novel latent variable model designed to capture text data generation process, as depicted on the left in Figure 1. The proposed generative model can be expressed probabilistically as:

$$p(\mathbf{x}, y) = \sum_{\mathbf{c}} p(\mathbf{x}|\mathbf{c}) p(y|\mathbf{c}) p(\mathbf{c}), \tag{1}$$

where we formulate human-interpretable concepts as latent variables c, and observed variables x and y represent text generated by the mapping g on c. We here distinguish observed variables into x and y, to align the input, i.e., context, and output token in next-token prediction framework, respectively.

Our goal is to consider realistic assumptions on the proposed latent variable model, ensuring the model's flexibility to closely approximate real text data. To this end, we do not impose any specific graph structure over latent variables. For example, we allow for any directed acyclic graph structure over latent variables from the perspective of causal representation learning. Most importantly, we highlight the following two factors that distinguish our model from existing latent variable models (Rajendran et al., 2024; Yan et al., 2024; Jin et al., 2020; Jiang et al., 2024), underscoring its novelty.

**Discrete Modeling for Text Data** Some existing latent variable models (Rajendran et al., 2024; Yan et al., 2024) for text data predominantly assume that both latent and observed variables are continuous, primarily to simplify identifiability analysis. However, this assumption may overlook the largely discrete characteristics of text data. In contrast, we assume throughout this work that all variables, both latent variables c and observed x and y, are discrete. This assumption can be justified for observed variables x and y, considering the inherently discrete structure of text data. Likewise, latent variables, which represent underlying concepts, can be well-suited to discrete modeling, as they often correspond to categorical distinctions or finite sets of semantic properties commonly observed in text data. For instance, in a topic modeling scenario, latent variables can represent distinct topics, such as "sports", "politics", or "technology", each of which forms a discrete category. This discrete structure aligns with the way humans typically interpret and classify information in text, making the assumption not only intuitive but also practical for real-world applications, as also assumed in prior work (Blei et al., 2003; Jin et al., 2020; Jiang et al., 2024).

No Invertibility Requirement Unlike existing approaches (Jiang et al., 2024; Rajendran et al., 2024), we do not require the mapping g from latent c to observed variables (x, y) to be invertible. Allowing non-invertibility is a deliberate design choice driven by fact that the mapping from latent to observed space is often complex and unknown. Imposing constraints on it may unnecessarily limit its flexibility. Moreover, in the context of text data, there are two key considerations. First, in natural language processing, the mapping from latent to observed space often involves a many-to-one relationship. For example, different combinations of emotional concepts can lead to the same sentiment label. In sentiment analysis, the combination of "positive sentiment" and "excitement" might result in an observed outcome such as "This is amazing!" (Miller, 1995). Second, some latent concepts may not be explicitly manifested in the observed sentence. For instance, a speaker's intent, tone, or implicit connotations may influence how a sentence is constructed, yet remain unobservable in surface-level text. This phenomenon is particularly prominent in pragmatics and discourse analysis, where unspoken contextual factors significantly shape meaning (Levinson, 1983). To formulate approximate invertibility, we define the mapping g from latent to observed space as follows:

**Definition 2.1.** We define the degree of approximate invertibility of the mapping  $\mathbf{g}$  from  $\mathbf{c}$  to  $(\mathbf{x}, y)$  by introducing an error term  $\epsilon$ , as follows:  $1 - p(\mathbf{c} = \mathbf{c}^* | \mathbf{x}, y) = \epsilon$ , where  $0 \le \epsilon < 1$ , and  $\mathbf{c}^* = \arg\max_{\mathbf{c}} p(\mathbf{c} | \mathbf{x}, y)$  represents the dominant mode of the posterior.

Remark 2.2. Such a relaxed condition naturally implies that we may not achieve exact identifiability results, as established in previous works, due to inevitable information loss in the context of non-invertible mappings. Nevertheless, this non-invertibility still allows for considering identifiability in an approximate sense, which will be discussed in Sec. 3.

# 3 THEORY: IDENTIFIABILITY OF THE PROPOSED LATENT VARIABLE MODEL

We now turn to the identifiability analysis of the proposed latent variable model, i.e., whether the latent variables  ${\bf c}$  can be uniquely recovered (up to an equivalence class) from observed data  ${\bf x}$  and  ${\bf y}$  under certain assumptions. We conduct this analysis within the next-token prediction framework, a widely adopted and empirically significant paradigm for training LLMs (Brown et al., 2020; Touvron et al., 2023b; Bi et al., 2024; Zhao et al., 2023).

We begin by introducing the general form of next-token prediction, which serves as the foundation for our analysis. The goal of next-token prediction is to learn a model that predicts the conditional distribution over the next token  $p(y|\mathbf{x})$ , which can be achieved by applying the softmax function over the logits produced by the model's final layer. This process mirrors multinomial logistic regression, where the model approximates the true conditional distribution  $p(y|\mathbf{x})$  by minimizing the crossentropy loss. In theory, assuming sufficient training data, a sufficiently expressive architecture and proper optimization, the model's predictions will converge to the true conditional distribution. We

emphasize that these assumptions are standard in the literature on identifiability analyses.

$$p(y|\mathbf{x}) = \frac{\exp\left(\mathbf{f}_{\mathbf{x}}(\mathbf{x})^T \mathbf{f}_{y}(y)\right)}{\sum_{y_j} \exp\left(\mathbf{f}_{\mathbf{x}}(\mathbf{x})^T \mathbf{f}_{y}(y_j)\right)},$$
 (2)

where  $y_j$  denotes a specific value of the observed variable y. The function  $\mathbf{f_x}(\mathbf{x})$  maps input  $\mathbf{x}$  to a representation space, and  $\mathbf{f}_y(y)$  corresponds to the model's final-layer weights, i.e., look-up table.

On the other hand, the true  $p(y|\mathbf{x})$  can be derived from Eq. 1 using Bayes' rule:

$$p(y|\mathbf{x}) = \sum_{\mathbf{c}} p(y|\mathbf{c})p(\mathbf{c}|\mathbf{x}). \tag{3}$$

By the expression for  $p(y|\mathbf{x})$  in both Eq. 2 and Eq. 3, we can align the right-hand sides of these two equations. Taking the logarithm of both sides yields:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x})^{T} \mathbf{f}_{y}(y) - \log \left( \sum_{y_{j}} \exp(\mathbf{f}_{\mathbf{x}}(\mathbf{x})^{T} \mathbf{f}_{y}(y_{j})) \right) = \log \sum_{\mathbf{c}} p(y|\mathbf{c}) p(\mathbf{c}|\mathbf{x}).$$
(4)

We now, as shown in Eq. 4, establish an initial connection between the LLM representations  $\mathbf{f}_{\mathbf{x}}$  in the inference model (left-hand side) and the latent variables  $\mathbf{c}$  in the generative model (right-hand side). To further study this connection, we introduce the following diversity condition.

**Diversity Condition** Let  $\mathbf{c} = [c^1, \dots, c^m]$  denote the latent variables, where  $c^k \in \mathcal{V}_k$  with  $|\mathcal{V}_k| = n_k, k = 1, \dots, m$ , so that  $\mathbf{c}$  can take  $\ell = \prod_{i=1}^m n_k$  possible values. Let  $\mathbf{c}_i, i = 1, \dots, \ell$ , denote all possible values of  $\mathbf{c}$ . We assume there exist  $\ell + 1$  distinct values of  $y, y_0, \dots, y_\ell$ , such that the following matrix  $\mathbf{L}$  is invertible, i.e.,

$$\mathbf{L} \in \mathbb{R}^{\ell \times \ell}, \qquad L_{j,i} = p(\mathbf{c} = \mathbf{c}_i \mid y = y_j) - p(\mathbf{c} = \mathbf{c}_i \mid y = y_0), \quad i, j = 1, \dots, \ell,$$

This assumption was originally developed in nonlinear independent component analysis (Hyvarinen & Morioka, 2016; Hyvarinen et al., 2019; Khemakhem et al., 2020). Intuitively, it requires that the conditional distribution  $p(\mathbf{c} \mid y)$  exhibits sufficiently diverse variation across different values of y. Similarly, we assume there exist  $\ell+1$  values of y such that the matrix

$$\hat{\mathbf{L}} = \left[ \mathbf{f}_y(y_1) - \mathbf{f}_y(y_0), \ \mathbf{f}_y(y_2) - \mathbf{f}_y(y_0), \ \dots, \ \mathbf{f}_y(y_\ell) - \mathbf{f}_y(y_0) \right] \in \mathbb{R}^{\ell \times \ell},$$

is invertible. This assumption has been employed in the context of identifiability analysis in inference space for LLMs (Roeder et al., 2021; Marconato et al., 2024). Justification can be found in Sec. C.

Under this diversity condition, we have the following identifiability result:

**Theorem 3.1.** Under the diversity condition above, the true latent variables c are related to LLM representations  $f_{\mathbf{x}}(\mathbf{x})$ , which are learned through the next-token prediction framework, by the following relationship:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x}) = \mathbf{A}[\log p(\mathbf{c} = \mathbf{c}_i | \mathbf{x})]_i + \mathbf{b} - (\hat{\mathbf{L}}^T)^{-1} \mathbf{h}_y,$$
 (5)

where  $\mathbf{h}_y = [h_{y_1} - h_{y_0}, ..., h_{y_\ell} - h_{y_0}]$  with  $h_{y_j} = [p(\mathbf{c} = \mathbf{c}_i | y = y_j)]_i^T [\log p(\mathbf{c} = \mathbf{c}_i | y = y_j, \mathbf{x})]_i$ , and  $\mathbf{b} = (\hat{\mathbf{L}}^T)^{-1}\mathbf{b}_y$ , with  $\mathbf{b}_y = [b(y = y_1) - b(y = y_0), ..., b(y = y_\ell) - b(y = y_0)]$  and  $b(y = y_j) = \mathbb{E}_{p(\mathbf{c}|y=y_j)}[\log p(y = y_j|\mathbf{c})]$ , and  $\mathbf{A} = (\hat{\mathbf{L}}^T)^{-1}\mathbf{L}$ . When  $\epsilon = 0$ ,  $\mathbf{f}_{\mathbf{x}}(\mathbf{x}) = \mathbf{A}[\log p(\mathbf{c} = \mathbf{c}_i | \mathbf{x})]_i + \mathbf{b}$ , and  $\epsilon \to 0$ ,  $\mathbf{f}_{\mathbf{x}}(\mathbf{x}) \approx \mathbf{A}[\log p(\mathbf{c} = \mathbf{c}_i | \mathbf{x})]_i + \mathbf{b}^{-1}$ .

Remark 3.2. This theorem establishes a precise relationship between the LLM representations  $\mathbf{f_x}(\mathbf{x})$  and the underlying latent concepts  $\mathbf{c}$ . This not only provides theoretical grounding for understanding LLM representations, but also offers a unified prospective for the linear representation hypothesis (Sec. 4). Beyond this, it suggests a principled approach for evaluating SAEs (Sec. 5).

#### 4 Theory $\rightarrow$ Insights: Unifying the Linear Hypothesis

In this section, we demonstrate how the identifiability result presented in Theorem 3.1 supports the linear representation hypothesis in LLMs. To this end, we first briefly introduce the linear representation hypothesis and then explain how it can be understood and unified through the linear matrix  $\mathbf{A}$  from our identifiability result.

<sup>&</sup>lt;sup>1</sup>Here,  $[\cdot]_i$  denotes a vector indexed by all latent configurations  $\mathbf{c}_i$ , and later we will also use  $[\cdot]_{c^k}$  to denote a vector indexed by all possible values of a single concept  $c^k$ .

# 4.1 THE LINEAR REPRESENTATION HYPOTHESIS

The linear representation hypothesis suggests that human-interpretable concepts in LLMs are represented linearly. This idea is supported by empirical evidence in various forms, including:

Concepts as Directions: Each concept is represented as a direction determined by the differences (i.e., vector offset) in representations of pairs that vary only in one latent concept of interest, e.g., gender. For instance, Rep("men")-Rep("women")  $\approx$  Rep("king")-Rep("queen") (Mikolov et al., 2013; Pennington et al., 2014; Turner et al., 2023).

Concept Manipulability: Previous studies have demonstrated that the value of a concept can be altered independently from others by introducing a corresponding steering vector (Li et al., 2024; Wang et al., 2023; Turner et al., 2023). For example, transitioning an output from a false to a truthful answer can be accomplished by adding a vector offset derived from counterfactual pairs that differ solely in the false/truthful concept.

**Linear Probing**: The value of a concept is often measured using a linear probe. For instance, the probability that the output language is French is logit-linear in the representation of the input. In this context, the linear weights can be interpreted as representing the concept of English/French (Park et al., 2023; Marks & Tegmark, 2023).

#### 4.2 Understanding the Linear Representation Hypothesis

The linear representation hypothesis has received growing empirical support in recent years. While recent work has aimed to develop unified frameworks for a deeper understanding of this phenomenon (Park et al., 2023; Marconato et al., 2024; Jiang et al., 2024), our approach seeks to explain it through the lens of identifiability. Before proceeding, we first provide the following definition:

**Definition 4.1.** We define the degree of approximate invertibility of the mapping from the latent variables  $\mathbf{c}$  to the observed variables  $\mathbf{x}$  by introducing an error term  $\epsilon_{\mathbf{x}}$ , as follows:  $1 - p(\mathbf{c} = \mathbf{c}_{\mathbf{x}}^*|\mathbf{x}) = \epsilon_{\mathbf{x}}$ , where  $0 < \epsilon_{\mathbf{x}} < 1$ , and  $\mathbf{c}_{\mathbf{x}}^* = \arg\max_{\mathbf{c}} p(\mathbf{c}|\mathbf{x})$  represents the dominant mode of  $p(\mathbf{c}|\mathbf{x})$ .

Next, we present two key corollaries derived from our identifiability results.

**Corollary 4.2** (Concepts Are Encoded in the Matrix **A**). Suppose that Theorem 3.1 holds, i.e.,  $\mathbf{f_x}(\mathbf{x}) \approx \mathbf{A} [\log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x})]_i + \mathbf{b}$ . Let  $\mathbf{x}_0$  and  $\mathbf{x}_1$  be a pair that differ only in the i-th concept  $c^k$ . Then as  $\epsilon_{\mathbf{x}} \to 0$ ,  $\mathbf{f_x}(\mathbf{x}_1) - \mathbf{f_x}(\mathbf{x}_0) \approx \tilde{\mathbf{A}}^k \left( [\log p(c^k \mid \mathbf{x}_1) - \log p(c^k \mid \mathbf{x}_0)]_{c^k} \right)$ , where  $\tilde{\mathbf{A}}^k = \mathbf{A}\mathbf{B}^k$ , where  $\mathbf{B}^k \in \{0,1\}^{\ell \times |\mathcal{V}_k|}$  is a binary lifting matrix that broadcasts each entry of  $[\log p(c^k \mid \mathbf{x})]_{c^k}$  to the corresponding index in  $[\log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x})]_i$ .

**Understanding Concepts as Directions** The corollary explains why the representation difference  $\operatorname{Rep}(\text{"man"}) - \operatorname{Rep}(\text{"woman"})$  can be closely approximated by  $\operatorname{Rep}(\text{"king"}) - \operatorname{Rep}(\text{"queen"})$ . In both the ("man", "woman") and ("king", "queen") pairs, the primary distinguishing factor is the latent concept of gender, while other concepts such as royalty remain largely unchanged. As a result, both  $\operatorname{Rep}(\text{"man"}) - \operatorname{Rep}(\text{"woman"})$  and  $\operatorname{Rep}(\text{"king"}) - \operatorname{Rep}(\text{"queen"})$  are driven by the same expression  $\tilde{\mathbf{A}}^k\left(\left[\log p(c^k\mid \mathbf{x}_1) - \log p(c^k\mid \mathbf{x}_0)\right]_{c^k}\right)$ . In the case of a binary concept  $c^k$ , the expression is further reduced to a vector direction corresponding to the concept of interest. See Appendix G.1 for details.

Understanding Concept Manipulability The corollary also supports the notion that a concept's value can be adjusted by adding a corresponding steering vector, such as Rep("man") – Rep("woman"). Adding the steering vector effectively modifies the original representation to produce a new representation, i.e.,  $\hat{\mathbf{f}}_{\mathbf{x}} = \mathbf{f}_{\mathbf{x}}(\mathbf{x}) + \alpha(\tilde{\mathbf{A}}^k([\log p(c^k \mid \mathbf{x}_1) - \log p(c^k \mid \mathbf{x}_0)]_{c^k})) = \mathbf{A}([\log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x})]_i + \alpha \mathbf{B}^k([\log p(c^k \mid \mathbf{x}_1) - \log p(c^k \mid \mathbf{x}_0)]_{c^k})) + \mathbf{b}$ , where  $\alpha$  represents an introduced weight (Wang et al., 2023; Turner et al., 2023). Thus, manipulating a concept via a steering vector is, in essence, equivalent to modifying the posterior distribution of the concept of interest given  $\mathbf{x}$ , which directly impacts the model's output.

To understand **Linear Probing**, we first introduce the following corollary:

**Corollary 4.3** (Linear Classifiability of Representations). Suppose that Theorem 3.1 holds, i.e.,  $\mathbf{f_x}(\mathbf{x}) \approx \mathbf{A} \left[ \log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x}) \right]_i + \mathbf{b}$ . Let  $\mathbf{x}_0$  and  $\mathbf{x}_1$  be pair data that differ only in the i-th concept

variable  $c^k$ . Then, as  $\epsilon_{\mathbf{x}} \to 0$ , the corresponding representations  $(\mathbf{f}_{\mathbf{x}}(\mathbf{x}_0), \mathbf{f}_{\mathbf{x}}(\mathbf{x}_1))$  are linearly separable along the variation of  $c^k$ . In particular, there exists a linear classifier with weight matrix  $\mathbf{W}$  such that  $\mathbf{W}\tilde{\mathbf{A}}^k \approx \mathbf{I}$ , and the corresponding logit is  $[p(c^k \mid \mathbf{x})]_{c^k}$ .

This corollary supports that latent concepts, e.g., English vs. French, can be reliably classified using a linear probing on the model's representations. This linear separability enables alignment between the model's predictive distribution and the true class distribution, achieved via cross-entropy minimization. A specific analysis for the binary case of concept  $c^k$  is provided in Appendix G.2.

#### 4.3 Unifying the Linear Representation Hypothesis

Taken together, Corollary 4.2 shows that the linear transformation  $\bf A$  underlies both concepts as directions and concept manipulability, while Corollary 4.3 demonstrates that linear probing is supported by a classifier  $\bf W$  satisfying  $\bf W \tilde{\bf A}^k \approx \bf I$ . This establishes a unified view on understanding the linear representation hypothesis: the various forms are all connected through the same underlying linear matrix  $\bf A$ , providing a unified theoretical perspective on how LLM representations linearly encode concepts. Importantly, according to Theorem 3.1, the matrix is defined as  $\bf A = (\hat{\bf L}^T)^{-1} \bf L$ , which, in essence, arises from the data diversity condition outlined in **Diversity Condition**.

# 5 Theory → Practice: Evaluating SAEs and Structured SAEs

**Evaluating SAEs** Broadly speaking, SAEs are designed with two primary objectives. First, they aim to learn a set of latent features  $\mathbf{z}$  such that sparse linear combinations  $\beta \mathbf{z}$  can accurately reconstruct LLM representations, i.e.,  $\mathbf{f_x}(\mathbf{x}) \approx \beta \mathbf{z}$ . Second, they seek to ensure that each learned feature  $z_i$  corresponds to an monosemantic, human-interpretable concept, thereby enabling a mechanistic understanding of LLMs. While reconstruction loss is commonly used to assess how well representations are reconstructed (Rajamanoharan et al., 2024a;b; Gao et al., 2025; Braun et al., 2024), it is a limited proxy for the second objective. A key challenge lies in the absence of ground truth for the underlying concepts (Kantamneni et al., 2025), making the evaluation of feature disentanglement challenging.

To address this issue, we propose a new evaluation method for SAEs grounded in our theoretical insights. Specifically, based on Theorem 3.1, the LLM representations  $\mathbf{f_x}(\mathbf{x})$  can be approximated as  $\mathbf{f_x}(\mathbf{x}) \approx \mathbf{A} \left[ \log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x}) \right]_i^2$ . Meanwhile, SAEs are trained to reconstruct  $\mathbf{f_x}(\mathbf{x})$  by  $\beta \mathbf{z}$ . Combining the two, we arrive at:  $\beta \mathbf{z} \approx \mathbf{A} \left[ \log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x}) \right]_i$ . This suggests that SAE features  $\mathbf{z}$  are linearly related to  $\left[ \log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x}) \right]_i$ . Therefore, if we expect each latent dimension  $z_i$  to encode only a single concept  $c^k$ , it should depend only on the posterior of that concept  $\log p(c^k \mid \mathbf{x})$ . Consequently, we can evaluate whether each  $z_i$  has successfully learned a monosemantic concept by measuring its linear correlation with  $\log p(c^k \mid \mathbf{x})$ . The question is: how can we obtain  $p(c^k \mid \mathbf{x})$ ?

Based on Corollary 4.3, this can be achieved using paired data that differ only in the k-th concept variable  $c^k$ . Specifically, we can construct paired data  $(\mathbf{x}_0, \mathbf{x}_1)$  that differ in only a single binary concept  $c^k$ , with labels  $c^k = 0$  for  $\mathbf{x}_0$  and  $c^k = 1$  for  $\mathbf{x}_1$ . We then train a linear classifier in a supervised manner on the corresponding LLM representations  $\mathbf{f}_{\mathbf{x}}(\mathbf{x}_0)$  and  $\mathbf{f}_{\mathbf{x}}(\mathbf{x}_1)$ , with the goal of predicting their labels. Once the classifier is trained, the resulting logit provides a estimate of the posterior probability  $p(c^k = 1 \mid \mathbf{x})$ . See Appendix I.2 for more rigorous details.

**Structured SAEs** Building on the evaluation method described above, our experiments (See Figure 4) suggest that sparsity regularization alone may may not be sufficient to fully disentangle the underlying concepts in LLM representations. Revisiting Theorem 3.1, the identifiability result depends on the proposed latent variable model. In this model, the complex dependencies in text data are encoded by the latent variables, whose interdependencies capture the underlying structure of the data. That is, latent variables are likely to exhibit strong interdependencies. Motivated by this, we propose exploring structured SAEs that incorporate additional regularization beyond sparsity to encourage learned features to model potential relationships among concepts.

In particular, we experiment with low-rank structures to complement sparsity, while noting that other forms of structured regularization could also be considered. Formally, structured SAEs minimize the

<sup>&</sup>lt;sup>2</sup>For simplicity, we omit the constant term **b** in this section, which does not affect the subsequent analysis.

following objective:

$$\mathcal{L} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{train}}} \left[ \|\mathbf{f}_{\mathbf{x}}(\mathbf{x}) - \overline{\mathbf{f}}_{\mathbf{x}}(\mathbf{x})\|_{2}^{2} + \lambda_{t} \left( \|\mathbf{S}\|_{p_{t}}^{p_{t}} + \gamma \|\mathbf{R}\|_{\text{nuc}} \right) \right], \tag{6}$$

where  $\bar{\mathbf{f}}_{\mathbf{x}}(\mathbf{x})$  denotes reconstruction,  $\lambda_t$  is the dynamically adjusted sparsity coefficient at step t, following the p-annealing SAE strategy (Karvonen et al., 2024),  $\gamma$  is a hyperparameter that balances the sparsity penalty and the low-rank regularization, the learned features  $\mathbf{z} = \mathbf{S} + \mathbf{R}$ ,  $\|\mathbf{S}\|_{p_t}^{p_t} = \sum_i |S_i|^{p_t}$  is the adaptive  $L_{p_t}$  norm promoting sparsity,  $\|\mathbf{R}\|_{\text{nuc}}$  is the nuclear norm, used to enforce a low-rank structure on  $\mathbf{R}$ . See Appendix I.1 for more implementation details.

# 6 EMPIRICAL EVALUATION ON SIMULATED AND REAL DATA WITH LLMS

**Simulation** We begin by conducting experiments on synthetic data, which is generated through the following process: First, we create random directed acyclic graphs (DAGs) with n latent variables, representing concepts. For each random DAG, the conditional probabilities of each variable given its parents are modeled using Bernoulli distributions, where the parameters are sampled uniformly from [0.2, 0.8]. To simulate a nonlinear mixture process, we then convert the latent variable samples into one-hot format and randomly apply a permutation matrix to the one-hot encoding, generating one-hot observed samples. These are then transformed into binary observed samples. To simulate next-token prediction, we randomly mask a part of the binary observed data, e.g.,,  $x_i$ , and predict it by use the remaining portion  $\mathbf{x}_{\setminus i}$ . Refer to Appendix H.1 for more details.

**Evaluation** In Theorem 3.1, we demonstrate that the representations learned by next-token prediction approximate a linear transformation of  $\log p(\mathbf{c}|\mathbf{x})$ , and this approximation becomes tighter when the mapping from  $\mathbf{c}$  to  $\mathbf{x}$  is approximately invertible. Building on this, Corollary 4.3 establishes that for a data pair  $(\mathbf{c}, \mathbf{x})$  differing only in the concept of interest, the representations  $\mathbf{f}_{\mathbf{x}}(\mathbf{x})$  is linearly separable. Therefore, to validate Theorem 3.1, we can assess the degree to which the learned representations can be classified linearly for data pairs  $(\mathbf{c}, \mathbf{x})$ .

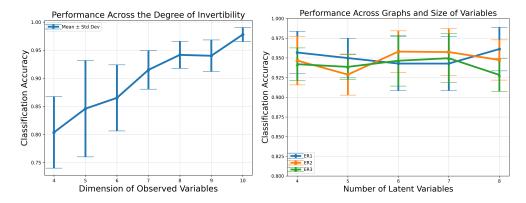


Figure 2: Classification accuracy vs. number of observed variables (left) and graph structures (right).

Our initial experiments investigate the relationship between the degree of invertibility of the mapping from  $\mathbf{c}$  to  $\mathbf{x}$  and the approximation of the identifiability result in Theorem 3.1. This exploration provides empirical insights into how invertibility influences the recovery of latent variables. To this end, we fix the size of the latent variables and gradually increase the size of the observed variables, thereby enhancing the degree of invertibility of the mapping from  $\mathbf{c}$  to  $\mathbf{x}$ . The left of Figure 2 demonstrates that classification accuracy improves as the size of the observed variables  $\mathbf{x}$  increases, aligning with results in Theorem 3.1.

We then examine the impact of latent graph structures on our identifiability results. To this end, we randomly generate DAG structures imposed on c. Specifically, random Erdős-Rényi (ER) graphs (ERDdS & R&wi, 1959) are generated with varying numbers of expected edges. For instance, ERk denotes graphs with d nodes and kd expected edges. The right panel of Figure 2 illustrates the relationship between classification accuracy and the size of the latent variables c under different settings, including ER1, ER2, and ER3. The results demonstrate that our identifiability findings

hold consistently across various graph structures and latent variable sizes, as evidenced by the linear classification accuracy.

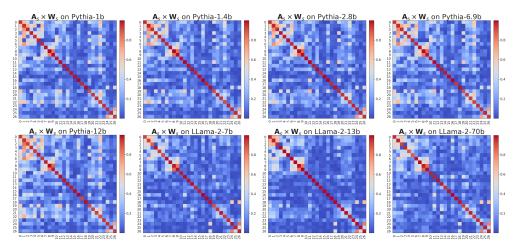


Figure 3: Results of the product  $\mathbf{A}_s \times \mathbf{W}_s$  across the LLaMA-2 and Pythia model families. Here,  $\mathbf{A}_s$  represents a matrix derived from the feature differences of 27 counterfactual pairs, while  $\mathbf{W}_s$  is a weight matrix obtained from a linear classifier trained on these features. The product approximates the identity matrix, supporting the theoretical findings in Corollary 4.3.

Experiments with LLMs We now present experiments on pre-trained LLMs. While it is challenging to collect all latent variables from real-world data to directly evaluate our identifiability result in Thorem 3.1, we can instead assess our corollaries to indirectly validate it. For Corollary 4.2, prior studies have already demonstrated the linear representation properties of LLM embeddings using counterfactual pairs that differ only in a single concept of interest (Mikolov et al., 2013; Pennington et al., 2014; Turner et al., 2023; Li et al., 2024; Wang et al., 2023; Park et al., 2023). Therefore, we shift our focus to a new property highlighted in Corollary 4.3, specifically the relationship between W and  $\tilde{\bf A}^k$ , which has not been explored in prior work. To investigate this, we require counterfactual pairs that differ only in a single concept. Constructing such counterfactual sentences is highly non-trivial, even for human annotators, due to the complexities of semantics and the need for precise control over contextual variations, as noted in prior studies (Park et al., 2023; Jiang et al., 2024). For our experiments, we utilize 27 counterfactual pairs from Park et al. (2023) that differ only in a binary concept, which are constructed based on the Big Analogy Test dataset (Gladkova et al., 2016). More details can be found in Sec. H.2 in Appendix.

Guided by Corollary 4.2, we use these 27 counterfactual pairs to obtain  $A^k$ . Specifically, for each pair differing in a binary concept, we compute the representation difference, and stacking all such vectors yields the matrix  $\mathbf{A}_s \in \mathbb{R}^{27 \times \text{dim}}$ , where each row corresponds to a concept direction, and dim denotes the representation dimension of the pre-trained LLM used. See a binary special case of Corollary 4.2 in Appendix G.1 for further details. To obtain W, motivated by Corollary 4.3, we train a linear classifier for each binary concept using the representations of the counterfactual pairs. The resulting weight vectors are stacked to form a matrix  $\mathbf{W}_s \in \mathbb{R}^{\dim \times 27}$ , where each column corresponds to the decision boundary for one concept. See a binary special case of Corollary 4.3 in Appendix G.2 for further details. Finally, we first normalize both  $A_s$  and  $W_s$  to remove the effect of scaling, which does not affect the semantics of a direction or decision boundary, and then examine the product  $A_s W_s$ . According to Corollary 4.3, the (i, j)-th entry of this product corresponds to the inner product between the representation difference vector for the i-th concept and the classifier weight vector for the j-th concept. When i = j, this inner product should be close to 1, indicating that the classifier is aligned with the concept direction. When  $i \neq j$ , the inner product should be less than 1, since the classifier for one concept should not be aligned with the concept direction of a different concept. Therefore, the matrix product  $A_sW_s$  should approximate the identity matrix I. Figure 3 displays the results of this product across the LLaMA-2 and Pythia model families. Refer to Appendix K for more results on LLaMA-3 and DeepSeek-R1. The results show that the product approximates the identity matrix, which is consistent with the theoretical finding in Corollary 4.3.

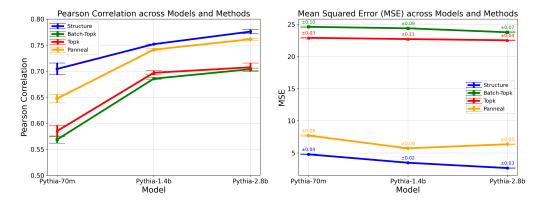


Figure 4: Comparison of SAE models. Pearson correlation scores (left) and MSE on the used counterfactual pairs (right) across different scales of Pythia.

**Experiments on SAEs** We finally conduct experiment on SAEs, training four sparse variants, including top-k SAE (Gao et al., 2025), batch-top-k SAE (Bussmann et al., 2024), p-annealing SAE (Karvonen et al., 2024), and the proposed structured SAE. Following Theorem 3.1, each SAE is trained on representations from the final hidden layer across Pythia-70m, Pythia-1.4b and Pythia-2.8b (Biderman et al., 2023), with *The Pile* corpus (Gao et al., 2020). For evaluation, as mentioned in Sec. 5, we use 27 counterfactual pairs from (Park et al., 2023) again to train a linear classifier using LogisticRegression implementation from the scikit-learn library, which provides logits, i.e., unnormalized  $p(c^k = 1|\mathbf{x})$ , for each of the 27 pairs. These counterfactual pairs are also passed through the trained SAEs to extract features  $\mathbf{z}$ . We search for the best-matching feature  $z_i$  for each  $p(c^k|\mathbf{x})$  according to Pearson correlation between  $\exp(z_i)$  and  $p(c^k|\mathbf{x})$ , to measure linear correlation.

Figure 4 reports the average Pearson correlations across the 27 counterfactual concepts (left, detailed results in Appendix I) and the reconstruction error measured by mean squared error (MSE, right). We draw two key conclusions. First, the proposed evaluation framework shows that: it differentiates SAE variants with sensitivity and provides interpretable evidence of how the learned features align with binary concepts. The trends are consistent with conventional reconstruction metrics such as MSE, reinforcing the reliability of our approach. Second, the results highlight the advantage of the proposed structured SAEs, which benefits from incorporating structured regularization. This advantage is consistently observed under both our evaluation framework and standard reconstruction metrics.

We emphasize again that obtaining counterfactual pairs is challenging. Even on this compact, high-precision benchmark, all four SAEs achieve Pearson correlations below 0.8 (1.0 indicates perfect recovery). Consequently, this benchmark is effective in differentiating the performance of the four SAEs. We hope it motivates the creation of more such high-quality counterfactual pairs.

#### 7 Conclusion

In this work, we propose a latent variable model to capture the generative process of text, representing high-level, human-interpretable concepts as latent variables. Under mild assumptions, our analysis provides a key insight into LLMs: a *linear property*, whereby LLM representations approximate a linear transformation of the posterior over latent variables. This establishes a foundational framework for understanding next-token prediction. Furthermore, we show that this linear property offers a unified theoretical perspective on various forms of the linear representation hypothesis, and motivates a principled evaluation strategy for sparse autoencoders. These findings open avenues for deeper exploration of how LLMs learn and represent complex patterns. Building on our results, we suggest the following future direction: develop methods to linearly unmix LLM representations to directly extract the probabilities of individual high-level concepts from the latent posterior. Further details are provided in Appendix L.

**Ethics Statement.** This research complies with the ICLR Code of Ethics. It does not use human subjects or sensitive personal data, and we are not aware of significant ethical risks.

**Reproducibility Statement.** We document all architectures, training procedures, and evaluation metrics. Code and scripts will be made publicly available to facilitate reproducibility after publication.

# REFERENCES

- Alberto Acerbi and Joseph M Stubbersfield. Large language models show human-like content biases in transmission chain experiments. *Proceedings of the National Academy of Sciences*, 120(44): e2313790120, 2023.
- Michal Aharon, Michael Elad, and Alfred Bruckstein. K-svd: An algorithm for designing over-complete dictionaries for sparse representation. *IEEE Transactions on signal processing*, 54(11): 4311–4322, 2006.
- Kartik Ahuja, Divyat Mahajan, Yixin Wang, and Yoshua Bengio. Interventional causal representation learning. In *ICML*, pp. 372–407. PMLR, 2023.
- Sanjeev Arora, Rong Ge, Tengyu Ma, and Ankur Moitra. Simple, efficient, and neural algorithms for sparse coding. In *Conference on learning theory*, pp. 113–149. PMLR, 2015.
- Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. A latent variable model approach to pmi-based word embeddings. *Transactions of the Association for Computational Linguistics*, 4:385–399, 2016.
- Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. Linear algebraic structure of word senses, with applications to polysemy. *Transactions of the Association for Computational Linguistics*, 6:483–495, 2018.
- Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*, 2024.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Dan Braun, Jordan Taylor, Nicholas Goldowsky-Dill, and Lee Sharkey. Identifying functionally important features with end-to-end sparse dictionary learning. *Advances in Neural Information Processing Systems*, 37:107286–107325, 2024.
- Johann Brehmer, Pim De Haan, Phillip Lippe, and Taco Cohen. Weakly supervised causal representation learning. *arXiv preprint arXiv:2203.16437*, 2022.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. https://transformer-circuits.pub/2023/monosemantic-features/index.html.

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
  - Simon Buchholz, Goutham Rajendran, Elan Rosenfeld, Bryon Aragam, Bernhard Schölkopf, and Pradeep Ravikumar. Learning linear causal representations from interventions under general nonlinear mixing. *arXiv preprint arXiv:2306.02235*, 2023.
  - Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. *arXiv preprint arXiv:2212.03827*, 2022.
  - Bart Bussmann, Patrick Leask, and Neel Nanda. Batchtopk sparse autoencoders. *arXiv preprint arXiv:2412.06410*, 2024.
  - Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
  - Bogdan Dumitrescu and Paul Irofti. *Dictionary learning algorithms and applications*. Springer, 2018.
  - Julian Eggert and Edgar Korner. Sparse coding and nmf. In 2004 IEEE International Joint Conference on Neural Networks (IEEE Cat. No. 04CH37541), volume 4, pp. 2529–2533. IEEE, 2004.
  - Michael Elad. Sparse and redundant representations: from theory to applications in signal and image processing. Springer Science & Business Media, 2010.
  - Michael Elad and Alfred M Bruckstein. A generalized uncertainty principle and sparse representation in pairs of bases. *IEEE Transactions on Information Theory*, 48(9):2558–2567, 2002.
  - Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022.
  - P ERDdS and A R&wi. On random graphs i. Publ. math. debrecen, 6(290-297):18, 1959.
  - Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
  - Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=tcsZt9ZNKD.
  - Anna Gladkova, Aleksandr Drozd, and Satoshi Matsuoka. Analogy-based detection of morphological and semantic relations with word embeddings: what works and what doesn't. In *Proceedings of the NAACL Student Research Workshop*, pp. 8–15, 2016.
  - Yuqi Gu and David B Dunson. Bayesian pyramids: Identifiable multilayer discrete latent structure models for discrete data. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 85(2):399–426, 2023.
  - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
  - Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610*, 2023.
  - Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. In *The Twelfth International Conference on Learning Representations*, 2024.

- Aapo Hyvarinen and Hiroshi Morioka. Unsupervised feature extraction by time-contrastive learning and nonlinear ica. *NeurIPS*, 29, 2016.
  - Aapo Hyvarinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ica using auxiliary variables and generalized contrastive learning. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 859–868. PMLR, 2019.
  - Yibo Jiang, Goutham Rajendran, Pradeep Ravikumar, Bryon Aragam, and Victor Veitch. On the origins of linear representations in large language models. *arXiv preprint arXiv:2403.03867*, 2024.
  - Shuning Jin, Sam Wiseman, Karl Stratos, and Karen Livescu. Discrete latent variable representations for low-resource text classification. *arXiv* preprint arXiv:2006.06226, 2020.
  - Subhash Kantamneni, Joshua Engels, Senthooran Rajamanoharan, Max Tegmark, and Neel Nanda. Are sparse autoencoders useful? a case study in sparse probing. *arXiv preprint arXiv:2502.16681*, 2025.
  - Adam Karvonen, Benjamin Wright, Can Rager, Rico Angell, Jannik Brinkmann, Logan Riggs Smith, Claudio Mayrink Verdun, David Bau, and Samuel Marks. Measuring progress in dictionary learning for language model interpretability with board game models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
  - Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoencoders and nonlinear ica: A unifying framework. In *AISTAS*, pp. 2207–2217. PMLR, 2020.
  - Bohdan Kivva, Goutham Rajendran, Pradeep Ravikumar, and Bryon Aragam. Learning latent causal graphs via mixture oracles. *Advances in Neural Information Processing Systems*, 34:18087–18101, 2021.
  - Lingjing Kong, Guangyi Chen, Biwei Huang, Eric P Xing, Yuejie Chi, and Kun Zhang. Learning discrete concepts in latent hierarchical models. *arXiv preprint arXiv:2406.00519*, 2024.
  - Stephen C Levinson. Pragmatics. Cambridge UP, 1983.
  - Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36, 2024.
  - Yuhang Liu, Zhen Zhang, Dong Gong, Mingming Gong, Biwei Huang, Anton van den Hengel, Kun Zhang, and Javen Qinfeng Shi. Identifying weight-variant latent causal models. *arXiv preprint arXiv:2208.14153*, 2022.
  - Yuhang Liu, Zhen Zhang, Dong Gong, Mingming Gong, Biwei Huang, Anton van den Hengel, Kun Zhang, and Javen Qinfeng Shi. Identifiable latent neural causal models. *arXiv preprint arXiv:2403.15711*, 2024a.
  - Yuhang Liu, Zhen Zhang, Dong Gong, Mingming Gong, Biwei Huang, Anton van den Hengel, Kun Zhang, and Javen Qinfeng Shi. Identifiable latent polynomial causal models through the lens of change. In *The Twelfth International Conference on Learning Representations*, 2024b.
  - Yuhang Liu, Zhen Zhang, Dong Gong, Biwei Huang, Mingming Gong, Anton van den Hengel, Kun Zhang, and Javen Qinfeng Shi. Revealing multimodal contrastive representation learning through latent partial causal models. *arXiv preprint arXiv:2402.06223*, 2024c.
  - Yuhang Liu, Zhen Zhang, Dong Gong, Mingming Gong, Biwei Huang, Anton van den Hengel, Kun Zhang, and Javen Qinfeng Shi. Latent covariate shift: Unlocking partial identifiability for multi-source domain adaptation. *Transactions on Machine Learning Research*, 2025.
  - Qing Lyu, Marianna Apidianaki, and Chris Callison-Burch. Representation of lexical stylistic features in language models' embedding space. *arXiv preprint arXiv:2305.18657*, 2023.
  - Christopher D Manning, Kevin Clark, John Hewitt, Urvashi Khandelwal, and Omer Levy. Emergent linguistic structure in artificial neural networks trained by self-supervision. *Proceedings of the National Academy of Sciences*, 117(48):30046–30054, 2020.

- Emanuele Marconato, Sébastien Lachapelle, Sebastian Weichwald, and Luigi Gresele. All or none: Identifiable linear properties of next-token predictors in language modeling. *arXiv preprint arXiv:2410.23501*, 2024.
  - Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. *arXiv preprint arXiv:2310.06824*, 2023.
  - Riccardo Massidda, Atticus Geiger, Thomas Icard, and Davide Bacciu. Causal abstraction with soft interventions. In *Conference on Causal Learning and Reasoning*, pp. 68–87. PMLR, 2023.
  - Tomáš Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 conference of the north american chapter of the association for computational linguistics: Human language technologies*, pp. 746–751, 2013.
  - George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11): 39–41, 1995.
  - Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello, and Emanuele Rodolà. Relative representations enable zero-shot latent space communication. *arXiv* preprint arXiv:2209.15430, 2022.
  - Neel Nanda, Andrew Lee, and Martin Wattenberg. Emergent linear representations in world models of self-supervised sequence models. *arXiv preprint arXiv:2309.00941*, 2023.
  - Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
  - Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. *arXiv* preprint arXiv:2311.03658, 2023.
  - Kiho Park, Yo Joong Choe, Yibo Jiang, and Victor Veitch. The geometry of categorical and hierarchical concepts in large language models. *arXiv preprint arXiv:2406.01506*, 2024.
  - Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532–1543, 2014.
  - Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Tom Lieberum, Vikrant Varma, János Kramár, Rohin Shah, and Neel Nanda. Improving dictionary learning with gated sparse autoencoders. *arXiv preprint arXiv:2404.16014*, 2024a.
  - Senthooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse autoencoders. *arXiv preprint arXiv:2407.14435*, 2024b.
  - Goutham Rajendran, Simon Buchholz, Bryon Aragam, Bernhard Schölkopf, and Pradeep Kumar Ravikumar. From causal to concept-based representation learning. In *NeurIPS*, 2024.
  - Geoffrey Roeder, Luke Metz, and Durk Kingma. On linear identifiability of learned representations. In *International Conference on Machine Learning*, pp. 9030–9039. PMLR, 2021.
  - Hassan Sajjad, Nadir Durrani, Fahim Dalvi, Firoj Alam, Abdul Rafae Khan, and Jia Xu. Analyzing encoded concepts in transformer language models. *arXiv preprint arXiv:2206.13289*, 2022.
  - Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
  - Anna Seigal, Chandler Squires, and Caroline Uhler. Linear causal disentanglement via interventions. *arXiv* preprint arXiv:2211.16467, 2022.
  - Lee Sharkey, Bilal Chughtai, Joshua Batson, Jack Lindsey, Jeff Wu, Lucius Bushnaq, Nicholas Goldowsky-Dill, Stefan Heimersheim, Alejandro Ortega, Joseph Bloom, et al. Open problems in mechanistic interpretability. *arXiv preprint arXiv:2501.16496*, 2025.

- Xinwei Shen, Furui Liu, Hanze Dong, Qing Lian, Zhitang Chen, and Tong Zhang. Weakly supervised disentangled generative causal representation learning. *The Journal of Machine Learning Research*, 23(1):10994–11048, 2022.
- Curt Tigges, Oskar John Hollinsworth, Atticus Geiger, and Neel Nanda. Linear representations of sentiment in large language models. *arXiv preprint arXiv:2310.15154*, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J Vazquez, Ulisse Mini, and Monte MacDiarmid. Activation addition: Steering language models without optimization. *arXiv e-prints*, pp. arXiv–2308, 2023.
- Burak Varici, Emre Acarturk, Karthikeyan Shanmugam, Abhishek Kumar, and Ali Tajer. Score-based causal representation learning with interventions. *arXiv preprint arXiv:2301.08230*, 2023.
- Julius Von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably isolates content from style. In *NeurIPS*, 2021.
- Julius von Kügelgen, Michel Besserve, Liang Wendong, Luigi Gresele, Armin Kekić, Elias Bareinboim, David Blei, and Bernhard Schölkopf. Nonparametric identifiability of causal representations from unknown interventions. *Advances in Neural Information Processing Systems*, 36, 2023.
- Zihao Wang, Lin Gui, Jeffrey Negrea, and Victor Veitch. Concept algebra for score-based conditional model. In *ICML 2023 Workshop on Structured Probabilistic Inference* {\&} *Generative Modeling*, 2023.
- Hanqi Yan, Lingjing Kong, Lin Gui, Yuejie Chi, Eric Xing, Yulan He, and Kun Zhang. Counterfactual generation with identifiability guarantees. *Advances in Neural Information Processing Systems*, 36, 2024.
- Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. Explainability for large language models: A survey. *ACM Transactions on Intelligent Systems and Technology*, 15(2):1–38, 2024.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv* preprint arXiv:2303.18223, 2023.
- Jieyu Zheng and Markus Meister. The unbearable slowness of being: Why do we live at 10 bits/s? *Neuron*, 2024.

# **Appendix**

# **Table of Contents**

A	Related Work	16
В	Limitations	17
C	Justification for the Diversity Condition	17
D	Proof of Theorem 3.1	18
E	Proof of Corollary 4.2	21
F	Proof of Corollary 4.3	23
G	Extension of Corollaries 4.2 and 4.3 for a Binary Concept G.1 Extension of Corollary 4.2 for a Binary Concept G.2 Extension of Corollary 4.3 for a Binary Concept	24 24 25
Н	Experimental Details Supporting Theoretical Results H.1 Simulation Details	26 26 26
I	Experiment on Sparse Autoencoders  I.1 Implementation of the Proposed Structured SAE  I.2 Details of Evaluation Metric	28 28 29 29
J	Further Discussion: Observations on LLMs and World Models  J.1 LLMs Mimic the Human World Model, Not the World Itself	31 31 32
	More Results on Llama-3 and DeepSeek-R1	34
L	Future Directions	35
M	Acknowledgment of LLMs Usage	35

# A RELATED WORK

810

811 812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830 831 832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853 854

855

856

857

858

859

860

861

862

863

Linearity of Representation in LLMs Recent studies have established the empirical finding that concepts in LLMs are often linearly encoded, a phenomenon known as the linear representation hypothesis (Mikolov et al., 2013; Pennington et al., 2014; Arora et al., 2016; Elhage et al., 2022; Burns et al., 2022; Tigges et al., 2023; Nanda et al., 2023; Moschella et al., 2022; Park et al., 2023; Li et al., 2024; Gurnee et al., 2023; Rajendran et al., 2024). Building on this observation, recent works (Park et al., 2023; Marconato et al., 2024) attempt to unify these findings into a cohesive framework, aiming to deepen our understanding and potentially inspire new insights. However, these works do not address why and how such linear properties emerge. Some previous works (Marconato et al., 2024; Roeder et al., 2021) have demonstrated identifiability results within the inference space, such as establishing connections between features derived from different inference models. However, these results are confined to the inference space and do not connect to the true latent variables in latent variable models. A recent work (Jiang et al., 2024) seeks to explain the origins of these linear properties but employs a latent variable model that differs from ours. From a technical perspective, their explanation is rooted in the implicit bias of gradient descent. In contrast, our work provides an explanation grounded in identifiability theory. This distinction highlights our focus on connecting the observed linear properties directly to the identifiability of true latent variables in latent variable models, offering a more comprehensive and theoretically robust understanding of these phenomena. In addition, recent work (Rajendran et al., 2024) assumes continuous latent and observed variables, while we assume both latent variables and observed variables to be discrete, aligning more closely with modeling natural language.

Causal Representation Learning This work is closely related to causal representation learning (Schölkopf et al., 2021), which aims to identify high-level latent causal variables from low-level observational data. Many prior studies (Brehmer et al., 2022; Von Kügelgen et al., 2021; Massidda et al., 2023; von Kügelgen et al., 2023; Ahuja et al., 2023; Seigal et al., 2022; Shen et al., 2022; Liu et al., 2022; Buchholz et al., 2023; Varici et al., 2023; Liu et al., 2024b; 2025; 2024a;c; Hyvarinen & Morioka, 2016; Hyvarinen et al., 2019; Khemakhem et al., 2020) have developed theoretical frameworks supporting the recovery of true latent variables up to simple transformations. However, these works primarily focus on continuous spaces and do not address the next-token prediction framework employed by LLMs. A subset of works (Gu & Dunson, 2023; Kong et al., 2024; Kivva et al., 2021) has explored causal representation learning in discrete spaces for both latent and observed variables, often imposing specific graph structures and assuming invertible mappings from latent to observed spaces. Despite their focus on discrete settings, none of these studies consider the next-token prediction framework in LLMs. In contrast, our work overcomes these limitations. We analyze approximate identifiability without relying on strict invertibility assumptions, which better aligns with the complex and often non-invertible relationships observed in real-world data. While a recent study has examined non-invertible mappings from latent to observed spaces, they rely on additional historical information to effectively restore invertibility. In our approach, we focus on achieving approximate identifiability without requiring such additional information. A very recent work (Rajendran et al., 2024) explores identifiability analysis for LLMs; however, they model both observed text data and latent variables as continuous and still assume invertibility in their analysis. In contrast, our framework explicitly considers discrete latent and observed variables, relaxed invertibility assumptions, and directly aligns with the next-token prediction paradigm, offering a more realistic and generalizable approach to causal representation learning in LLMs.

**Sparse Autoencoders** Polysemanticity, a phenomenon observed in recent studies, roughly speaking, refers to cases where a single representation encodes multiple distinct, human-interpretable concepts (Arora et al., 2018). Early investigations suggested that neural networks represent features by linear superposition and motivated efforts to disentangle human-interpretable concepts from such linear mixing (Elhage et al., 2022). This can be achieved by using sparse autoencoders (Huben et al., 2024; Rajamanoharan et al., 2024a; Gao et al., 2025; Braun et al., 2024; Bricken et al., 2023), a technique closely related to the well-known framework of dictionary learning (Dumitrescu & Irofti, 2018; Eggert & Korner, 2004; Elad, 2010; Elad & Bruckstein, 2002; Aharon et al., 2006; Arora et al., 2015). In contrast to these works, we propose structured SAEs to model the dependencies among latent concepts. Furthermore, motivated by our theoretical findings, we introduce a new evaluation method for SAEs, grounded in our justified theoretical results.

# **B** LIMITATIONS

 One of the main contributions of this work is establishing an identifiability result for the next-token prediction framework, a widely used approach for training LLMs. This result hinges on a key assumption: the diversity condition, which requires the data distribution to exhibit sufficiently strong and diverse variations. The diversity condition was originally introduced in the context of nonlinear ICA (Hyvarinen & Morioka, 2016; Hyvarinen et al., 2019; Khemakhem et al., 2020) and has since been adopted in the identifiability analysis within the causal representation learning community. In essence, it has emerged as a fundamental requirement for identifiability. This condition might be satisfied given that the training data for LLMs is primarily sourced from the Internet, encompassing a broad and diverse range of content. Recent empirical analyses of LLMs further support the plausibility of this assumption (Roeder et al., 2021; Marconato et al., 2024), reinforcing its relevance in the context of identifiability.

Similar to existing works (Park et al., 2023; Jiang et al., 2024; Marconato et al., 2024; Rajendran et al., 2024), our analysis is limited to the last layer in LLMs and does not provide a justification for intermediate layers. Formalizing the relationship between representations in intermediate layers and identifiability analysis presents additional challenges. This complexity arises from the intricate nature of intermediate layer representations. Therefore, extending identifiability results to intermediate layers remains an open and more challenging problem.

#### C JUSTIFICATION FOR THE DIVERSITY CONDITION

The diversity condition involves two key requirements: (1) the invertibility of  $\mathbf{L}$ , which pertains to the latent generative process, and (2) the invertibility of  $\hat{\mathbf{L}}$ , which pertains to the latent inference process. The first condition was originally developed in the context of nonlinear independent component analysis (Hyvarinen & Morioka, 2016; Hyvarinen et al., 2019; Khemakhem et al., 2020) and has been further extended to identifiability analysis in the causal representation learning community (Liu et al., 2022; 2024b). Intuitively, the invertibility of  $\mathbf{L}$  ensures that the latent variables  $\mathbf{c}$  exhibit sufficiently strong and diverse variation across different labels y, making it possible, in principle, to recover  $\mathbf{c}$  uniquely from y. Furthermore, the invertibility of  $\hat{\mathbf{L}}$  indicates that there exists a sufficiently diverse set of y values such that the corresponding difference vectors  $\mathbf{f}_y(y_j) - \mathbf{f}_y(y_0)$  form a linearly independent set spanning the image of  $\mathbf{f}_y$ . This assumption is generally mild: as noted by Roeder et al. (2021), the probability that randomly initialized and stochastically updated parameters of  $\mathbf{f}_y$  produce linearly dependent difference vectors is effectively zero.

We emphasize that, while the diversity condition is primarily a theoretical assumption ensuring identifiability, its plausibility can be assessed empirically using pre-trained LLMs on real data. Directly observing all latent variables c is infeasible, but the corollaries derived from our theory, including Corollary 4.2 and Corollary 4.3, provide predictions that can be indirectly validated. Prior work has shown that LLM embeddings exhibit linear representation properties when comparing counterfactual pairs differing in a single concept (Mikolov et al., 2013; Pennington et al., 2014; Turner et al., 2023; Li et al., 2024; Wang et al., 2023; Park et al., 2023), supporting the theoretical prediction of sufficient diversity in latent directions (Corollary 4.2). In addition, our experiments show that the alignment between representation differences and classifier weights across multiple LLM families (Figures 3 and 6) closely matches the predictions of Corollary 4.3, suggesting that the diversity condition might plausibly hold in practice.

# D Proof of Theorem 3.1

*Proof.* Next-token prediction mirrors multinomial logistic regression, where the model approximates the true conditional distribution, i.e.,  $p(y|\mathbf{x})$ , by minimizing the cross-entropy loss. In this framework, given sufficient training data, a sufficiently expressive architecture, and effective optimization, the model's predictions are expected to converge to the true conditional distribution, as follows:

$$p(y|\mathbf{x}) = \frac{\exp(\mathbf{f}_{\mathbf{x}}(\mathbf{x})^T \mathbf{f}_y(y))}{\sum_{y'} \exp(\mathbf{f}_{\mathbf{x}}(\mathbf{x})^T \mathbf{f}_y(y_i))}.$$
 (7)

On the other hand, for the proposed latent variable model described in Eq. 1), the true conditional distribution  $p(y|\mathbf{x})$  can be derived using Bayes' rule:

$$p(y|\mathbf{x}) = \sum_{\mathbf{c}} p(y|\mathbf{c})p(\mathbf{c}|\mathbf{x}). \tag{8}$$

By comparing Eq. 8 and Eq. 7, we arrive at the following relationship:

$$\frac{\exp(\mathbf{f}_{\mathbf{x}}(\mathbf{x})^T \mathbf{f}_y(y))}{\sum_{y'} \exp(\mathbf{f}_{\mathbf{x}}(\mathbf{x})^T \mathbf{f}_y(y_y))} = \sum_{\mathbf{c}} p(y|\mathbf{c}) p(\mathbf{c}|\mathbf{x}).$$
(9)

Taking the logarithm on both sides of Eq. 9, we obtain:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x})^{T}\mathbf{f}_{y}(y) - \log Z(\mathbf{x}) = \log \left( \sum_{c} p(y|\mathbf{c})p(\mathbf{c}|\mathbf{x}) \right)$$
(10)

where  $Z(\mathbf{x}) = \sum_{y_j} \exp(\mathbf{f}_{\mathbf{x}}(\mathbf{x})^T \mathbf{f}_y(y_j))$  represents the normalization constant. This transformation provides a direct link between the representations  $\mathbf{f}(\mathbf{x})$  learned by next-token prediction (Eq. 7) and the true posterior distribution  $p(\mathbf{c}|\mathbf{x})$ .

Now, let us focus on the right of Eq. 10, we can obtain that:

$$\log\left(\sum_{c} p(y|\mathbf{c})p(\mathbf{c}|\mathbf{x})\right) = \log p(y|\mathbf{x}),\tag{11}$$

$$= \mathbb{E}_{p(\mathbf{c}|y)} [\log p(y|\mathbf{x})], \tag{12}$$

$$= \mathbb{E}_{p(\mathbf{c}|y)} \left[ \log \left( p(\mathbf{c}|y, \mathbf{x}) \frac{p(y|\mathbf{x})}{p(\mathbf{c}|y, \mathbf{x})} \right) \right], \tag{13}$$

$$= \mathbb{E}_{p(\mathbf{c}|y)} \left[ \log p(y, \mathbf{c}|\mathbf{x}) \right] - \mathbb{E}_{p(\mathbf{c}|y)} \left[ \log p(\mathbf{c}|y, \mathbf{x}) \right], \tag{14}$$

$$= \mathbb{E}_{p(\mathbf{c}|y)} \left[ \log p(\mathbf{c}|\mathbf{x}) \right] + \mathbb{E}_{p(\mathbf{c}|y)} \left[ \log \underbrace{p(y|\mathbf{c},\mathbf{x})}_{p(y|\mathbf{c})} \right] - \mathbb{E}_{p(\mathbf{c}|y)} \left[ \log p(\mathbf{c}|y,\mathbf{x}) \right],$$

where  $p(y|\mathbf{c}, \mathbf{x}) = p(y|\mathbf{c})$ , due to the conditional independence of y and  $\mathbf{x}$  given  $\mathbf{c}$ , as defined in the proposed latent variable model.

Together with Eq. 15 and Eq. 10, we have:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x})^{T}\mathbf{f}_{y}(y) - \log Z(\mathbf{x}) = \mathbb{E}_{p(\mathbf{c}|y)}[\log p(\mathbf{c}|\mathbf{x})] + \underbrace{\mathbb{E}_{p(\mathbf{c}|y)}[\log p(y|\mathbf{c})]}_{b_{y}} - \mathbb{E}_{p(\mathbf{c}|y)}[\log p(\mathbf{c}|y,\mathbf{x})].$$
(16)

Define  $\mathbb{E}_{p(\mathbf{c}|y)}[\log p(y|\mathbf{c})] = b_y$ , and define a vector  $[p(\mathbf{c} = \mathbf{c}_i|\mathbf{x})]_i$  as the vector constructed by the probabilities of all possible values of  $\mathbf{c}$ , conditional on  $\mathbf{x}$ . As a result,  $\mathbb{E}_{p(\mathbf{c}|y)}\log p(\mathbf{c}|\mathbf{x}) = \sum_{\mathbf{c}} p(\mathbf{c}|y)\log p(\mathbf{c}|\mathbf{x}) = [p(\mathbf{c} = \mathbf{c}_i|y)]_i^T[\log p(\mathbf{c} = \mathbf{c}_i|\mathbf{x})]_i$ , and similarly  $\mathbb{E}_{p(\mathbf{c}|y)}\log p(\mathbf{c}|y,\mathbf{x}) = [p(\mathbf{c} = \mathbf{c}_i|y)]_i^T[\log p(\mathbf{c} = \mathbf{c}_i|y,\mathbf{x})]_i$ . Then we can re-write Eq. 16 as follows:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x})^{T} \mathbf{f}_{y}(y) - \log Z(\mathbf{x}) = [p(\mathbf{c} = \mathbf{c}_{i}|y)]_{i}^{T} [\log p(\mathbf{c} = \mathbf{c}_{i}|\mathbf{x})]_{i}$$
(17)

$$-[p(\mathbf{c} = \mathbf{c}_i|y)]_i^T[\log p(\mathbf{c} = \mathbf{c}_i|y,\mathbf{x})]_i + b_y$$
(18)

(15)

Let  $y_0, \ldots, y_\ell$  be the points provided by the diversity condition outlined in Section 3, for y = 0, we have:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x})^{T}\mathbf{f}_{y}(y_{0}) - \log Z(\mathbf{x}) = [p(\mathbf{c} = \mathbf{c}_{i}|y = y_{0})]_{i}^{T}[\log p(\mathbf{c} = \mathbf{c}_{i}|\mathbf{x})]_{i} - \underbrace{[p(\mathbf{c} = \mathbf{c}_{i}|y = y_{0})]_{i}^{T}[\log p(\mathbf{c} = \mathbf{c}_{i}|y = y_{0}, \mathbf{x})]_{i}}_{h_{y_{0}}} + b_{y_{0}},$$
(19)

where we define  $h_{y_0} = [p(\mathbf{c} = \mathbf{c}_i | y = y_0)]_i^T [\log p(\mathbf{c} = \mathbf{c}_i | y = y_0, \mathbf{x})]_i$ . For y = 1, we similarly obtain:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x})^{T} \mathbf{f}_{y}(y_{1}) - \log Z(\mathbf{x}) = [p(\mathbf{c} = \mathbf{c}_{i} | y = y_{1})]_{i}^{T} [\log p(\mathbf{c} = \mathbf{c}_{i} | \mathbf{x})]_{i}$$
$$-h_{y_{1}} + h_{y_{1}}. \tag{20}$$

Subtracting Eq. 19 from Eq. 20, we get the following expression:

$$(\mathbf{f}_y(y_1)^T - \mathbf{f}_y(y_0)^T) \mathbf{f}_\mathbf{x}(\mathbf{x}) = ([p(\mathbf{c} = \mathbf{c}_i | y = y_1) - p(\mathbf{c} = \mathbf{c}_i | y = y_0)]_i^T) [\log p(\mathbf{c} = \mathbf{c}_i | \mathbf{x})]_i$$

$$- (h_{y_1} - h_{y_0}) + b_{y_1} - b_{y_0}.$$

$$(21)$$

According to the diversity condition, where y can take  $\ell+1$  values, we can obtain a total of  $\ell$  equations similar to Eq. 21. Collecting all of these equations, we have:

$$\underbrace{\left(\mathbf{f}_{y}(y_{1}) - \mathbf{f}_{y}(y_{0}), \dots, \mathbf{f}_{y}(y_{\ell}) - \mathbf{f}_{y}(y_{0})\right)^{T}}_{\hat{\mathbf{L}}^{T}} \mathbf{f}_{\mathbf{x}}(\mathbf{x}) \qquad (22)$$

$$= \underbrace{\left(\left[p(\mathbf{c} = \mathbf{c}_{i}|y = y_{1}) - p(\mathbf{c} = \mathbf{c}_{i}|y = y_{0})\right]_{i}, \dots, \left[p(\mathbf{c} = \mathbf{c}_{i}|y = y_{\ell}) - p(\mathbf{c} = \mathbf{c}_{i}|y = y_{0})\right]_{i}\right)^{T}}_{\mathbf{L}} \times \left[\log p(\mathbf{c} = \mathbf{c}_{i}|\mathbf{x})\right]_{i} - \underbrace{\left[h_{y_{1}} - h_{y_{0}}, \dots, h_{y_{\ell}} - h_{y_{0}}\right]}_{\mathbf{h}_{y}} + \underbrace{\left[b_{y_{1}} - b_{y_{0}}, \dots, b_{y_{\ell}} - b_{y_{0}}\right]}_{\mathbf{b}_{y}}. \qquad (23)$$

According to the diversity condition, the matrix  $\hat{\mathbf{L}}$  of size  $\ell \times \ell$  is invertible, as a result, we arrive:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x}) = \underbrace{(\hat{\mathbf{L}}^T)^{-1}\mathbf{L}}_{\mathbf{A}}[\log p(\mathbf{c} = \mathbf{c}_i|\mathbf{x})]_i - (\hat{\mathbf{L}}^T)^{-1}\mathbf{h}_y + (\hat{\mathbf{L}}^T)^{-1}\mathbf{b}_y.$$
(24)

Here due to the diversity condition, the matrix  $\mathbf{L} = (p(\mathbf{c}|y=y_1) - p(\mathbf{c}|y=y_0), ..., p(\mathbf{c}|y=y_\ell) - p(\mathbf{c}|y=y_0))$  of size  $\ell \times \ell$  is also invertible,  $\mathbf{A}$  is invertible.

Now, we focus on the term  $\mathbf{b}_y$  on the right-hand side of Eq. 24. Note that  $b_{y_i} = \mathbb{E}_{p(\mathbf{c}|y_i)}[\cdot]$  is a constant with respect to  $\mathbf{c}$ , as the expectation integrates over all possible values of  $\mathbf{c}$ . As a result, the entire term  $(\hat{\mathbf{L}}^T)^{-1}\mathbf{b}_y$ , denoted as  $\mathbf{b}$ , is also a constant.

We next examine the term  $\mathbf{h}_y$  in Eq. 24, considering the cases where the mapping from the latent space to the observed space is exactly invertible and where it is only approximately invertible.

**Invertible** According to definition 2.1, when the mapping g from latent space to observed space is invertible, meaning that for

$$1 - p(\mathbf{c} = \mathbf{c}^* \mid \mathbf{x}, y) = \epsilon, \tag{25}$$

we have  $\epsilon = 0$ . Then, for  $\mathbf{h}_{u}$ , we analyze each component, i.e.,  $h_{u_i} - h_{u_0}$ , where

$$h_{y_i} = \mathbb{E}_{p(\mathbf{c}|y=y_i)} \left[ \log p(\mathbf{c} \mid \mathbf{x}, y=y_i) \right], \quad h_{y_0} = \mathbb{E}_{p(\mathbf{c}|y=y_0)} \left[ \log p(\mathbf{c} \mid \mathbf{x}, y=y_0) \right]. \tag{26}$$

When  $\epsilon = 0$ , the posterior distribution  $p(\mathbf{c} \mid \mathbf{x}, y)$  becomes a delta distribution centered at  $\mathbf{c}^*$ , i.e.,

$$p(\mathbf{c} \mid \mathbf{x}, y) = \delta(\mathbf{c} - \mathbf{c}^*), \tag{27}$$

which implies that the posterior is concentrated at a single point  $\mathbf{c}^*$ , satisfying  $\log p(\mathbf{c}^* \mid \mathbf{x}, y) = 0$ . In this case, we have

$$h_{u_i} - h_{u_0} = 0. (28)$$

**Approximately Invertible** When the mapping from latent space to observed space is approximately invertible, i.e.,  $\epsilon \to 0$ , for  $\mathbf{h}_{u}$ , we analyze the difference:

$$h_{y_i} - h_{y_0} = \mathbb{E}_{p(\mathbf{c}|y=y_i)} \left[ \log p(\mathbf{c} \mid \mathbf{x}, y_i) \right] - \mathbb{E}_{p(\mathbf{c}|y=y_0)} \left[ \log p(\mathbf{c} \mid \mathbf{x}, y_0) \right]. \tag{29}$$

Since the generative map  $\mathbf{g}: \mathbf{c} \mapsto (\mathbf{x}, y)$  is approximately invertible, the posterior  $p(\mathbf{c} \mid \mathbf{x}, y)$  becomes sharply concentrated at a unique  $\mathbf{c}^*$ .

Then for each expectation:

$$\mathbb{E}_{p(\mathbf{c}|y)}\left[\log p(\mathbf{c} \mid \mathbf{x}, y)\right] = p(\mathbf{c}^* \mid y)\log(1 - \epsilon) + \sum_{\mathbf{c} \neq \mathbf{c}^*} p(\mathbf{c} \mid y)\log p(\mathbf{c} \mid \mathbf{x}, y). \tag{30}$$

Now note:

- $\log(1 \epsilon) \approx 0$ , when  $\epsilon \to 0$ ,
- $\log p(\mathbf{c} \mid \mathbf{x}, y) \leq \log \epsilon \text{ for all } \mathbf{c} \neq \mathbf{c}^*,$
- and for  $\mathbf{c} \neq \mathbf{c}^*$ ,  $p(\mathbf{c} \mid y)$  is bounded.

Thus, when  $\epsilon \to 0$ , the entire expectation satisfies:

$$\mathbb{E}_{p(\mathbf{c}|y)}\left[\log p(\mathbf{c}\mid \mathbf{x}, y)\right] \approx \log \epsilon,\tag{31}$$

which do not depend on y. Therefore,

$$h_{y_i} - h_{y_0} \xrightarrow[\epsilon \to 0]{} 0. \tag{32}$$

# E Proof of Corollary 4.2

*Proof.* We first prove that: when  $\epsilon_{\mathbf{x}} \to 0$ , i.e.,  $p(\mathbf{c} \mid \mathbf{x})$  becomes sharply peaked at  $\mathbf{c}_{\mathbf{x}}^*$ , we can approximate:

$$p(\mathbf{c}) \approx p(c^k \mid \mathbf{x}) \cdot p(\mathbf{c}^{-k} \mid \mathbf{x}),$$
 (33)

where  $\mathbf{c}^{-k}$  denotes all concepts except  $c^k$ . To this end, we analysis their KL (Kullback–Leibler) residual term:

Residual := 
$$D_{KL} \left( p(\mathbf{c} \mid \mathbf{x}) \parallel p(c^k \mid \mathbf{x}) p(\mathbf{c}^{-k} \mid \mathbf{x}) \right)$$
. (34)

Since when  $\epsilon_{\mathbf{x}} \to 0$ ,  $p(\mathbf{c} \mid \mathbf{x})$  is sharply peaked at a particular configuration  $\mathbf{c}_{x}^{*} = (c^{k*}, \mathbf{c}^{-k*})$ , i.e.,:

$$p(\mathbf{c} \mid \mathbf{x}) = \begin{cases} 1 - \epsilon_{\mathbf{x}}, & \text{if } \mathbf{c} = \mathbf{c}_{x}^{*}, \\ \epsilon_{\mathbf{x}} \cdot r(\mathbf{c}), & \text{otherwise,} \end{cases}$$
(35)

where  $\sum_{\mathbf{c}\neq\mathbf{c}^*} r(\mathbf{c}) = 1$ .

**Main Term in KL** The KL divergence is given by:

$$D_{KL}(p(\mathbf{c} \mid \mathbf{x}) || p(c^k \mid \mathbf{x}) p(\mathbf{c}^{-k} \mid \mathbf{x})) = \sum_{\mathbf{c}} p(\mathbf{c} \mid \mathbf{x}) \log \frac{p(\mathbf{c} \mid \mathbf{x})}{p(c^k \mid \mathbf{x}) p(\mathbf{c}^{-k} \mid \mathbf{x})}.$$
 (36)

The main contribution comes from  $\mathbf{c} = \mathbf{c}_x^*$ :

$$(1 - \epsilon_{\mathbf{x}}) \log \frac{1 - \epsilon_{\mathbf{x}}}{p(c^{k*} \mid \mathbf{x}) \cdot p(\mathbf{c}^{-k*} \mid \mathbf{x})}.$$
 (37)

Note that:

$$p(c^{k*} \mid \mathbf{x}) = \sum_{\mathbf{c}^{-k}} p(c^{k*}, \mathbf{c}^{-k} \mid \mathbf{x}) \ge p(c^{k*}, \mathbf{c}^{-k*} \mid \mathbf{x}) = 1 - \epsilon_{\mathbf{x}},$$
 (38)

and similarly:

$$p(\mathbf{c}^{-k*} \mid \mathbf{x}) \ge 1 - \epsilon_{\mathbf{x}}.\tag{39}$$

Hence:

$$p(c^{k*} \mid \mathbf{x}) \cdot p(\mathbf{c}^{-k*} \mid \mathbf{x}) \ge (1 - \epsilon_{\mathbf{x}})^2, \tag{40}$$

$$\Rightarrow \frac{1 - \epsilon_{\mathbf{x}}}{p(c^{k*} \mid \mathbf{x}) \cdot p(\mathbf{c}^{-k*} \mid \mathbf{x})} \le \frac{1}{1 - \epsilon_{\mathbf{x}}}.$$
 (41)

Thus, the main term becomes:

$$(1 - \epsilon_{\mathbf{x}}) \log \frac{1 - \epsilon_{\mathbf{x}}}{p(c^{k*} \mid \mathbf{x}) \cdot p(\mathbf{c}^{-k*} \mid \mathbf{x})} \le (1 - \epsilon_{\mathbf{x}}) \log \left(\frac{1}{1 - \epsilon_{\mathbf{x}}}\right) = -(1 - \epsilon_{\mathbf{x}}) \log(1 - \epsilon_{\mathbf{x}}). \tag{42}$$

**Tail Term** For  $\mathbf{c} \neq \mathbf{c}_x^*$ , the contribution to the KL divergence is:

$$\epsilon_{\mathbf{x}} r(\mathbf{c}) \log \left( \frac{\epsilon_{\mathbf{x}} r(\mathbf{c})}{\epsilon_{\mathbf{x}}^2 r(c^k) r(\mathbf{c}^{-k})} \right) = \epsilon_{\mathbf{x}} r(\mathbf{c}) \log \left( \frac{1}{\epsilon_{\mathbf{x}}} \cdot \frac{r(\mathbf{c})}{r(c^k) r(\mathbf{c}^{-k})} \right). \tag{43}$$

Summing over all  $\mathbf{c} \neq \mathbf{c}_{x}^{*}$ , we obtain:

$$\epsilon_{\mathbf{x}} \sum_{\mathbf{c} \neq \mathbf{c}_{x}^{*}} r(\mathbf{c}) \log \left( \frac{1}{\epsilon_{\mathbf{x}}} \cdot \frac{r(\mathbf{c})}{r(c^{k})r(\mathbf{c}^{-k})} \right) = \epsilon_{\mathbf{x}} \log \frac{1}{\epsilon_{\mathbf{x}}} \sum_{\mathbf{c} \neq \mathbf{c}_{x}^{*}} r(\mathbf{c}) + \epsilon_{\mathbf{x}} \sum_{\mathbf{c} \neq \mathbf{c}_{x}^{*}} r(\mathbf{c}) \log \left( \frac{r(\mathbf{c})}{r(c^{k})r(\mathbf{c}^{-k})} \right).$$
(44)

Since the term  $r(\mathbf{c})$  is bounded, we conclude:

Tail Term = 
$$o(\epsilon_{\mathbf{x}} \log \epsilon_{\mathbf{x}})$$
. (45)

which vanishes faster than the main term as  $\epsilon_{\mathbf{x}} \to 0$ .

Final Bound Combining both contributions in Main Term and Tail Term, we get:

$$D_{\mathrm{KL}}\left(p(\mathbf{c} \mid \mathbf{x}) \mid p(c^{k} \mid \mathbf{x}) \cdot p(\mathbf{c}^{-k} \mid \mathbf{x})\right) \le -(1 - \epsilon_{\mathbf{x}}) \log(1 - \epsilon_{\mathbf{x}}) + o(\epsilon_{\mathbf{x}} \log \epsilon_{\mathbf{x}}). \tag{46}$$

Here when  $\epsilon_{\mathbf{x}} \to 0$ ,  $D_{\mathrm{KL}}\left(p(\mathbf{c} \mid \mathbf{x}) \parallel p(c^k \mid \mathbf{x}) \cdot p(\mathbf{c}^{-k} \mid \mathbf{x})\right) \to 0$ .

Now that we have shown that Eq. 33 holds, we can rewrite the term  $[\log p(\mathbf{c}_i \mid \mathbf{x})]_i$  as:

$$\left[\log p(\mathbf{c}_i \mid \mathbf{x})\right]_i \approx \mathbf{B}^k \left[\log p(c^k \mid \mathbf{x})\right]_{c^k} + \mathbf{B}^{-i} \left[\log p(\mathbf{c}^{-k} \mid \mathbf{x})\right]_{\mathbf{c}^{-k}}.$$
 (47)

Here,  $\mathbf{B}^k$  and  $\mathbf{B}^{-k}$  are binary broadcasting matrices that expand the marginal log-probability vectors  $\left[\log p(\mathbf{c}^k\mid\mathbf{x})\right]_{\mathbf{c}^k}$  and  $\left[\log p(\mathbf{c}^{-k}\mid\mathbf{x})\right]_{\mathbf{c}^{-k}}$  to the full configuration space  $\left[\log p(\mathbf{c}_i\mid\mathbf{x})\right]_i$ , respectively.

As a result, for pair  $\mathbf{x}_0$  and  $\mathbf{x}_1$  that differ only in the *i*-th concept variable  $c^k$ , their difference in posterior distribution is:

$$[\log p(\mathbf{c}_i \mid \mathbf{x}_1) - \log p(\mathbf{c}_i \mid \mathbf{x}_0)]_i$$

$$\approx \mathbf{B}^k \left[\log p(c^k \mid \mathbf{x}_1) - \log p(c^k \mid \mathbf{x}_0)\right]_{c^k} + \mathbf{B}^{-k} \left[\log p(\mathbf{c}^{-k} \mid \mathbf{x}_1) - \log p(\mathbf{c}^{-k} \mid \mathbf{x}_0)\right]_{\mathbf{c}^{-k}}. \quad (48)$$

We now show that:

$$p(\mathbf{c}^{-k} \mid \mathbf{x}_1) - p(\mathbf{c}^{-k} \mid \mathbf{x}_0) \to 0, \quad \text{as } \epsilon_{\mathbf{x}} \to 0$$
 (49)

so the term  $\mathbf{B}^{-i} \left[ \log p(\mathbf{c}^{-i} \mid \mathbf{x}_1) - \log p(\mathbf{c}^{-i} \mid \mathbf{x}_0) \right]_{\mathbf{c}^{-i}}$  in Eq. 48 vanishes.

Recall Eq. 35, we have

$$p(\mathbf{c} \mid \mathbf{x}) = \begin{cases} 1 - \epsilon_{\mathbf{x}}, & \text{if } \mathbf{c} = \mathbf{c}_{x}^{*} = (c^{k*}, \mathbf{c}^{-k*}), \\ \epsilon_{\mathbf{x}} \cdot r(\mathbf{c}), & \text{otherwise}, \end{cases}$$
(50)

Then, for  $\mathbf{c}^{-i}$ , we have:

$$p(\mathbf{c}^{-k} \mid \mathbf{x}) = \sum_{c^k} p(c^k, \mathbf{c}^{-k} \mid \mathbf{x}). \tag{51}$$

Combining this with Eq. 50, we have:

$$p(\mathbf{c}^{-k*} \mid \mathbf{x}) = 1 - \epsilon_{\mathbf{x}} + \epsilon_{\mathbf{x}} \sum_{c^k \neq c^{k*}} r(c^k, \mathbf{c}^{-k*}) = 1 - \epsilon_{\mathbf{x}} + o(\epsilon_{\mathbf{x}}),$$
 (52)

$$p(\mathbf{c}^{-k} \mid \mathbf{x}) = \epsilon_{\mathbf{x}} \sum_{c^k} r(c^k, \mathbf{c}^{-k}) = o(\epsilon_{\mathbf{x}}), \quad \text{for } \mathbf{c}^{-k} \neq \mathbf{c}^{-k*}.$$
(53)

Both Eq. 52 and Eq. 52 only depends on  $\epsilon$ . As a result, if  $\mathbf{x}_0$  and  $\mathbf{x}_1$  have the same values on the components relevant to  $\mathbf{c}^{-k}$ , the difference between  $p(\mathbf{c}^{-k} \mid \mathbf{x}_1)$  and  $p(\mathbf{c}^{-k} \mid \mathbf{x}_0)$  is of order  $o(\epsilon_{\mathbf{x}})$ , as:

$$p(\mathbf{c}^{-i} \mid \mathbf{x}_1) = p(\mathbf{c}^{-i} \mid \mathbf{x}_0) + o(\epsilon_{\mathbf{x}}), \tag{54}$$

which implies:

$$\log p(\mathbf{c}^{-k} \mid \mathbf{x}_1) - \log p(\mathbf{c}^{-k} \mid \mathbf{x}_0) = \log \left( 1 + \frac{o(\epsilon_{\mathbf{x}})}{p(\mathbf{c}^{-k} \mid \mathbf{x}_0)} \right) = o(\epsilon_{\mathbf{x}}). \tag{55}$$

Consequently,

$$\mathbf{B}^{-k} \left[ \log p(\mathbf{c}^{-k} \mid \mathbf{x}_1) - \log p(\mathbf{c}^{-k} \mid \mathbf{x}_0) \right]_{\mathbf{c}^{-k}} = o(\epsilon_{\mathbf{x}}) \to 0 \quad \text{as } \epsilon_{\mathbf{x}} \to 0.$$
 (56)

Together with Eq. 48 and the result from Theorem 3.1, we get:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x}_1) - \mathbf{f}_{\mathbf{x}}(\mathbf{x}_0) \approx \mathbf{A}\mathbf{B}^k \left( \left[ \log p(c^k \mid \mathbf{x}_1) \right]_{c^k} - \left[ \log p(c^k \mid \mathbf{x}_0) \right]_{c^k} \right). \tag{57}$$

PROOF OF COROLLARY 4.3 *Proof.* Again, when  $\epsilon_{\mathbf{x}} \to 0$ , i.e., when  $p(\mathbf{c} \mid \mathbf{x})$  becomes sharply peaked at  $\mathbf{c}_{\mathbf{x}}^*$ , we can approximate:  $p(\mathbf{c} \mid \mathbf{x}) \approx p(c^k \mid \mathbf{x}) \cdot p(\mathbf{c}^{-k} \mid \mathbf{x}),$ (58)Given the above, neglecting the constant term in the result in Theorem 3.1,  $\mathbf{f}_{\mathbf{x}}(\mathbf{x}) \approx \mathbf{A} \left[ \log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x}) \right]_i$ (59)which can be rewritten as  $\mathbf{f}_{\mathbf{x}}(\mathbf{x}) \approx \mathbf{A} \left[ \log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x}) \right]_i \approx \mathbf{A}(\mathbf{B}^k \left[ \log p(c^k \mid \mathbf{x}) \right]_{c^k} + \mathbf{B}^{-k} \left[ \log p(\mathbf{c}^{-k} \mid \mathbf{x}) \right]_{c^{-k}}).$ (60)In this case, for a data pair  $(\mathbf{x}_0, \mathbf{x}_1)$  that differ only in the latent variable  $c^k$ , the representations  $\mathbf{f}_{\mathbf{x}}(\mathbf{x})$ are passed to a linear classifier with weights W. The classifier produces the logits: logits  $\approx \mathbf{W} (\mathbf{A} (\mathbf{B}^k \lceil \log p(c^k \mid \mathbf{x}) \rceil_{c^k} + \mathbf{B}^{-k} \lceil \log p(\mathbf{c}^{-k} \mid \mathbf{x}) \rceil_{\mathbf{c}^{-k}})).$ (61)For correct classification under cross entropy loss, the logits must match the true probabilities:  $\mathbf{logits} = \left[ p(c^k \mid \mathbf{x}) \right]_{c^k},$ (62)where we omit constant scaling factors for simplicity, corresponding to the normalization applied prior to the softmax operation in the cross-entropy loss. Combining Eq. 61 and Eq. 62, the weight matrix W must satisfy the condition:  $\mathbf{W}(\mathbf{A}\mathbf{B}^k) \approx \mathbf{I},$ (63)which ensures that the classifier produces the correct logits. Here  $\mathbf{I}$  denotes the identify matrix.  $\Box$ 

# G EXTENSION OF COROLLARIES 4.2 AND 4.3 FOR A BINARY CONCEPT

When considering pair that differ only in a concept of interest, i.e.,  $c^k$ , and assuming that the concept is binary, both Corollaries 4.2 and 4.3 can be further refined, yielding the following results.

#### G.1 EXTENSION OF COROLLARY 4.2 FOR A BINARY CONCEPT

**Corollary G.1** (Binary Concept Direction). Suppose that Theorem 3.1 holds, and let  $c^k$  be a binary concept variable, i.e.,  $c^k \in \{0,1\}$ . Let  $\mathbf{x}_0$  and  $\mathbf{x}_1$  be a pair of inputs that differ only in the i-th binary concept  $c^k$ , with  $c^k = 0$  for  $\mathbf{x}_0$  and  $c^k = 1$  for  $\mathbf{x}_1$ . Then, as  $\epsilon_{\mathbf{x}} \to 0$ , the representation difference simplifies as:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x}_{1}) - \mathbf{f}_{\mathbf{x}}(\mathbf{x}_{0}) \approx \tilde{\mathbf{A}}^{k} \left( \left[ \log p(c^{k} \mid \mathbf{x}_{1}) - \log p(c^{k} \mid \mathbf{x}_{0}) \right]_{c^{k}} \right) \approx \log p(c^{k} = 0 \mid \mathbf{x}_{1}) \cdot \tilde{\mathbf{A}}^{k} \begin{bmatrix} 1 \\ -1 \end{bmatrix},$$

$$or, \approx \log p(c^{k} = 1 \mid \mathbf{x}_{0}) \cdot \tilde{\mathbf{A}}^{k} \begin{bmatrix} 1 \\ -1 \end{bmatrix},$$

$$(64)$$

where  $\tilde{\mathbf{A}}^k = \mathbf{A}\mathbf{B}^k$ ,  $\mathbf{B}^k$  is a binary lifting matrix that broadcasts each entry of  $[\log p(c^k \mid \mathbf{x})]_{c^k}$  to the corresponding index in  $[\log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x})]_i$ . This shows that changes in a binary concept are encoded in a specific direction in the representation space defined by  $\tilde{\mathbf{A}}^k$ .

*Proof.* Recall Eq. 35, the joint concept distribution conditioned on input x follows a peaked structure:

$$p(\mathbf{c} \mid \mathbf{x}) = \begin{cases} 1 - \epsilon_{\mathbf{x}}, & \text{if } \mathbf{c} = \mathbf{c}_{\mathbf{x}}^*, \\ \epsilon_{\mathbf{x}} \cdot r(\mathbf{c}), & \text{otherwise,} \end{cases}$$
 (66)

where  $\mathbf{c}_{\mathbf{x}}^* = (c^{k*}, \mathbf{c}^{-k*})$  is the dominant concept configuration for  $\mathbf{x}$ , and  $r(\mathbf{c})$  is a normalized residual distribution over all non-dominant  $\mathbf{c}$ .

Let  $\mathbf{x}_0$  and  $\mathbf{x}_1$  differ only in the *i*-th binary concept variable  $c^k$ , with:

$$\mathbf{c}_{\mathbf{x}_0}^* = (0, \mathbf{c}^{-k*}), \quad \mathbf{c}_{\mathbf{x}_1}^* = (1, \mathbf{c}^{-k*}).$$
 (67)

Then the marginal probabilities for  $c^k$  are:

For  $\mathbf{x}_1$ :

$$p(c^k = 1 \mid \mathbf{x}_1) = (1 - \epsilon_{\mathbf{x}_1}) + \epsilon_{\mathbf{x}_1} \cdot \sum_{\mathbf{c}^{-k} \neq \mathbf{c}^{-k*}} r(1, \mathbf{c}^{-i}) = 1 - \epsilon_{\mathbf{x}_1} \cdot \alpha_1, \tag{68}$$

where  $\alpha_1 := 1 - \sum_{\mathbf{c}^{-k} \neq \mathbf{c}^{-k*}} r(1, \mathbf{c}^{-k})$ .

$$p(c^k = 0 \mid \mathbf{x}_1) = 1 - p(c^k = 1 \mid \mathbf{x}_1) = \epsilon_{\mathbf{x}_1} \cdot \alpha_1.$$
 (69)

For  $\mathbf{x}_0$ :

$$p(c^k = 0 \mid \mathbf{x}_0) = (1 - \epsilon_{\mathbf{x}_0}) + \epsilon_{\mathbf{x}_0} \cdot \sum_{\mathbf{c}^{-k} \neq \mathbf{c}^{-k*}} r(0, \mathbf{c}^{-k}) = 1 - \epsilon_{\mathbf{x}_0} \cdot \alpha_0, \tag{70}$$

where  $\alpha_0 := 1 - \sum_{\mathbf{c}^{-k} \neq \mathbf{c}^{-k*}} r(0, \mathbf{c}^{-k})$ .

$$p(c^k = 1 \mid \mathbf{x}_0) = 1 - p(c^k = 1 \mid \mathbf{x}_0) = \epsilon_{\mathbf{x}_0} \cdot \alpha_0.$$
 (71)

Taking logarithmic, we have:

$$\log p(c^k = 1 \mid \mathbf{x}_1) = \log(1 - \epsilon_{\mathbf{x}_1} \cdot \alpha_1) \approx 0, \quad \text{as } \epsilon_{\mathbf{x}} \to 0$$
 (72)

$$\log p(c^k = 0 \mid \mathbf{x}_1) = \log(\epsilon_{\mathbf{x}_1} \cdot \alpha_1). \tag{73}$$

$$\log p(c^k = 0 \mid \mathbf{x}_0) = \log(1 - \epsilon_{\mathbf{x}_0} \cdot \alpha_0) \approx 0, \quad \text{as } \epsilon_{\mathbf{x}} \to 0$$
 (74)

$$\log p(c^k = 1 \mid \mathbf{x}_0) = \log(\epsilon_{\mathbf{x}_0} \cdot \alpha_0) \approx \log(\epsilon_{\mathbf{x}_1} \cdot \alpha_1) = \log p(c^k = 0 \mid \mathbf{x}_1), \quad \text{as } \epsilon_{\mathbf{x}} \to 0.$$
 (75)

Then the vector difference:

$$\left[\log p(c^{k} \mid \mathbf{x}_{1}) - \log p(c^{k} \mid \mathbf{x}_{0})\right]_{c^{k}} = \left[\log p(c^{k} = 0 \mid \mathbf{x}_{1}) - \log p(c^{k} = 0 \mid \mathbf{x}_{0})\right]$$

$$\log p(c^{k} \mid \mathbf{x}_{1}) - \log p(c^{k} = 1 \mid \mathbf{x}_{0})\right]$$
(76)

$$\approx \begin{bmatrix} \log p(c^k = 0 \mid \mathbf{x}_1) \\ -\log p(c^k = 1 \mid \mathbf{x}_0) \end{bmatrix}$$
 (77)

$$\approx \log p(c^k = 0 \mid \mathbf{x}_1) \begin{bmatrix} 1 \\ -1 \end{bmatrix} \tag{78}$$

Finally, the representation difference becomes:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x}_1) - \mathbf{f}_{\mathbf{x}}(\mathbf{x}_0) \approx \tilde{\mathbf{A}}^k \left( \left[ \log p(c^k \mid \mathbf{x}_1) - \log p(c^k \mid \mathbf{x}_0) \right]_{c^k} \right) \approx \log p(c^k = 0 \mid \mathbf{x}_1) \cdot \tilde{\mathbf{A}}^k \begin{bmatrix} 1 \\ -1 \end{bmatrix}. \tag{79}$$

Here, note that: 
$$\log p(c^k = 0 \mid \mathbf{x}_1) \approx \log p(c^k = 1 \mid \mathbf{x}_0)$$
 as shown in Eq. 75.

# G.2 EXTENSION OF COROLLARY 4.3 FOR A BINARY CONCEPT

**Corollary G.2** (Binary Concept Classification). Suppose that Theorem 3.1 holds, i.e.,  $\mathbf{f_x}(\mathbf{x}) \approx \mathbf{A} \left[ \log p(\mathbf{c} = \mathbf{c}_i \mid \mathbf{x}) \right]_i + \mathbf{b}$ . Let  $\mathbf{x}_0$  and  $\mathbf{x}_1$  be pair data that differ only in the i-th binary concept variable  $c^k$ , with labels  $c^k$ , where  $c^k = 0$  for  $\mathbf{x}_0$  and  $c^k = 1$  for  $\mathbf{x}_1$ . Then when  $\epsilon_{\mathbf{x}} \to 0$ , the corresponding representations  $(\mathbf{f}(\mathbf{x}_0), \mathbf{f}(\mathbf{x}_1))$  are linearly separable with a weight vector  $\mathbf{w}$  satisfying  $\mathbf{w}^{\top} \tilde{\mathbf{A}}^k \begin{bmatrix} -1 \\ 1 \end{bmatrix} \approx 1$ . The corresponding logit is the (unnormalized)  $p(c^k = 1 \mid \mathbf{x})$ .

*Proof.* When  $\epsilon_{\mathbf{x}} \to 0$ , the posterior  $p(\mathbf{c} \mid \mathbf{x})$  becomes sharply peaked at a unique mode  $\mathbf{c}_{\mathbf{x}}^*$ . This implies a near-independence of  $c^k$  and  $\mathbf{c}^{-i}$  given  $\mathbf{x}$ , allowing us to write:

$$p(\mathbf{c} \mid \mathbf{x}) \approx p(c^k \mid \mathbf{x}) \cdot p(\mathbf{c}^{-k} \mid \mathbf{x}).$$
 (80)

Taking logs and substituting into Theorem 3.1 (neglecting the bias term b), we obtain:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x}) \approx \mathbf{A} \left[ \log p(\mathbf{c}_i \mid \mathbf{x}) \right]_i \approx \mathbf{A} \left( \mathbf{B}^k \left[ \log p(c^k \mid \mathbf{x}) \right]_{c^k} + \mathbf{B}^{-i} \left[ \log p(\mathbf{c}^{-k} \mid \mathbf{x}) \right]_{\mathbf{c}^{-k}} \right),$$
 (81)

where  $\mathbf{B}^k$  and  $\mathbf{B}^{-k}$  denote the lifting operators that map the marginal log-probabilities of  $c^k$  and  $\mathbf{c}^{-k}$  into the joint log-probability vector space.

Define  $\tilde{\mathbf{A}}^k := \mathbf{A}\mathbf{B}^k$ . Then:

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x}) \approx \tilde{\mathbf{A}}^k [\log p(c^k \mid \mathbf{x})]_{c^k} + \text{(terms involving only } \mathbf{c}^{-k}).$$
 (82)

Consider a linear classifier with weight vector  $\mathbf{w}$  applied to  $\mathbf{f}_{\mathbf{x}}(\mathbf{x})$ :

logit = 
$$\mathbf{w} \cdot \mathbf{f}_{\mathbf{x}}(\mathbf{x}) \approx \mathbf{w} \tilde{\mathbf{A}}^k \begin{bmatrix} p(c^k = 0 \mid \mathbf{x}) \\ 1 - p(c^k = 0 \mid \mathbf{x}) \end{bmatrix} + \text{const.}$$
 (83)

Let  $\mathbf{w}^{\top} \tilde{\mathbf{A}}^k = [s_0, s_1]$ . Then:

logit 
$$\approx s_0 \log p(c^k = 0 \mid \mathbf{x}) + s_1 \log p(c^k = 1 \mid \mathbf{x}) + \text{const} = (s_1 - s_0) \log p(c^k = 1 \mid \mathbf{x}) + \text{const.}$$
(84)

For the classifier to correctly separate  $x_0$  and  $x_1$  under cross-entropy loss, we require:

$$logit = log p(c^k = 1 \mid \mathbf{x}), \tag{85}$$

where we omit the normalization constant prior to the softmax, since it does not affect the analysis. This is equivalent to:

$$\mathbf{w}^{\top} \tilde{\mathbf{A}}^k \begin{bmatrix} -1\\1 \end{bmatrix} \approx 1, \tag{86}$$

This completes the proof.  $\Box$ 

# H EXPERIMENTAL DETAILS SUPPORTING THEORETICAL RESULTS

#### H.1 SIMULATION DETAILS

For the left side of Figure 2, which investigates the relationship between the degree of invertibility in the mapping from c to x and the approximation of the identifiability result in Theorem 3.1, we aim to exclude other uncertain factors that might affect the result. To achieve this, we keep the number of latent variables constant (i.e., 3) and ensure that the graph structure follows a chain structure. Based on this structure, we model the conditional probabilities of each variable, given its parents, using Bernoulli distributions. The parameters of these distributions are uniformly sampled from the interval [0.2, 0.8], which are then used to generate the latent variables. Subsequently, we apply a one-hot encoding to these samples to obtain one-hot formal representations. These one-hot samples are then randomly permuted. For 3 latent variables, there are  $2^3$  possible permutations, each corresponding to 3 observed binary variables, resulting in a total of  $2^3 \times 3$  different observed binary variables. We then randomly sample from these observed variables, varying the sample size. For example, as shown on the left in Figure 2, we can select different variables as observed variables. Clearly, as the number of observed variables increases, the mutual information between observed variables increases.

For the right side of Figure 2, we explore the robustness of our identifiability result in Theorem 3.1 with respect to both the graph structure and the size of the latent variables. To this end, we randomly generate DAG structures in the latent space using Erdős-Rényi (ER) graphs (ERDdS & R&wi, 1959), where ERk denotes graphs with d nodes and kd expected edges. For each ERk configuration, we also vary the size of the latent variables from 4 to 8, allowing us to examine how the size of latent variables influences the identifiability results. In terms of the observed variables, we adapt the experimental setup from the left side of Figure 2 to determine the appropriate observed variable size for different latent variable sizes. This ensures that the degree of invertibility from the latent space to the observed space remains sufficiently high, a crucial factor for the accuracy of our identifiability analysis.

Throughout the simulation, we use the following: In each experiment, we randomly mask one observed variable  $x_i$ , and use the remaining observed variables to predict it. Specifically, the remaining variables and the corresponding mask matrix are used as inputs to an embedding layer. This embedding layer transforms the input into a high-dimensional feature representation. The generated embeddings are then passed through a Multi-Layer Perceptron (MLP)-based architecture to extract meaningful features, e.g.,  $\mathbf{f_x}(\mathbf{x})$ . The MLP model consists of three layers, each with 256 hidden units. After ach layer, we apply Batch Normalization unit to stabilize training and a ReLU nonlinear activation function to introduce nonlinearity. The final output of the MLP is used to predict the masked variable through a linear classification layer. This allows us to assess how well the model can predict missing or masked values based on the remaining observed variables. We employ the Adam optimizer with a learning rate of 1e-4. To ensure robustness and account for potential variability in the results, we conduct each experimental setting with five different runs, each initialized with a different random seed. This procedure helps mitigate the effects of random initialization and provides a more reliable evaluation of the model's performance.

For evaluation, we use the LogisticRegression classifier from the scikit-learn library, which operates on the features extracted from the output of the MLP-based architecture described above.

# H.2 EXPERIMENTAL DETAILS ON LLMS

Unlike in simulation studies, where we have access to the complete set of latent variables, in real-world scenarios, their true values remain inherently unknown. This limitation arises from the nature of latent variables, they are unobserved and must be inferred indirectly from the data. As a consequence, we cannot directly validate the linear identifiability results established in Theorem 3.1, since such validation would require explicit knowledge of these latent variables.

However, we can instead verify Corollary 4.3, which is a direct consequence of Theorem 3.1. By doing so, we provide indirect empirical evidence supporting the theoretical identifiability results. To achieve this, we need collect counterfactual pairs of data instances that differ in controlled and specific ways. These counterfactual pairs are essential for testing the implications of our theory in the context of real-world data.

Generating such counterfactual pairs, however, presents significant challenges. First, the inherent complexity and nuances of natural language make it difficult to create pairs that differ in precisely the intended contexts while leaving other aspects unchanged. Second, as highlighted in prior works (Park et al., 2023; Jiang et al., 2024), constructing such counterfactual sentences is a highly non-trivial task, even for human annotators, due to the intricacies of semantics and the need for precise control over contextual variations.

We utilize the 27 counterfactual pairs introduced in (Park et al., 2023), which provide a structured and well-curated set of counterfactual pairs. These concepts encompass a wide range of semantic and morphological transformations, as detailed in Table 1. By leveraging this established dataset, we ensure consistency with previous research while facilitating a robust and meaningful evaluation of Corollary 4.3.

We first use these 27 counterfactual pairs to construct a  $A_s$  with size  $27 \times dim$  by using the differences in the representations of these 27 counterfactual pairs, where dim corresponds to the feature dimension of the used LLM, such as 4096 for the Llama-27B model. To construct the corresponding matrix  $W_s$ , we train a linear classifier using these 27 counterfactual pairs. Specifically, we use the representations of the counterfactual pairs as input and the corresponding values of the latent variables as output. As a result, the corresponding linear weights for the 27 counterfactual pairs are be used to create  $W_s$ . In our experiments, we employ the LogisticRegression classifier from the scikit-learn library.

Note that, since the 27 counterfactual pairs involve only binary concepts, each row of  $\mathbf{A}_s$  corresponds to the direction of a concept associated with a counterfactual pair, as stated in Corollary G.1. Similarly, as discussed in Corollary G.2, each column of  $\mathbf{W}_s$  denotes the classifier direction for a specific pair and aligns with the corresponding row of  $\mathbf{A}_s$  for that concept. As a result, after normalizing  $\mathbf{A}_s$  and  $\mathbf{W}_s$  to remove the arbitrary logit scaling (which is irrelevant before the softmax in cross-entropy loss), we expect to observe that  $\mathbf{A}_s\mathbf{W}_s\approx\mathbf{I}$ . Therefore, we apply normalization to both  $\mathbf{A}_s$  and  $\mathbf{W}_s$  prior to computing their product.

Table 1: Concept names, one example of the counterfactual pairs, and the number of used pairs, taken from (Park et al., 2023).

#	Concept	Example	Word Pair Counts
1	$verb \Rightarrow 3pSg$	(accept, accepts)	50
2	verb ⇒ Ving	(add, adding)	50
3	$verb \Rightarrow Ved$	(accept, accepted)	50
4	$Ving \Rightarrow 3pSg$	(adding, adds)	50
5	$Ving \Rightarrow Ved$	(adding, added)	50
6	$3pSg \Rightarrow Ved$	(adds, added)	50
7	$verb \Rightarrow V + able$	(accept, acceptable)	50
8	$verb \Rightarrow V + er$	(begin, beginner)	50
9	$verb \Rightarrow V + tion$	(compile, compilation)	50
10	$verb \Rightarrow V + ment$	(agree, agreement)	50
11	$adj \Rightarrow un + adj$	(able, unable)	50
12	$adj \Rightarrow adj + ly$	(according, accordingly)	50
13	small ⇒ big	(brief, long)	25
14	$thing \Rightarrow color$	(ant, black)	50
15	$thing \Rightarrow part$	(bus, seats)	50
16	$country \Rightarrow capital$	(Austria, Vienna)	158
17	$pronoun \Rightarrow possessive$	(he, his)	4
18	$male \Rightarrow female$	(actor, actress)	52
19	lower ⇒ upper	(always, Always)	73
20	noun ⇒ plural	(album, albums)	100
21	$adj \Rightarrow comparative$	(bad, worse)	87
22	$adj \Rightarrow superlative$	(bad, worst)	87
23	$frequent \Rightarrow infrequent$	(bad, terrible)	86
24	$English \Rightarrow French$	(April, avril)	116
25	French ⇒ German	(ami, Freund)	128
26	French $\Rightarrow$ Spanish	(annee, año)	180
27	$German \Rightarrow Spanish$	(Arbeit, trabajo)	228

# I EXPERIMENT ON SPARSE AUTOENCODERS

# I.1 IMPLEMENTATION OF THE PROPOSED STRUCTURED SAE

The proposed structured SAE employs two regularization terms: a structured regularization to model the dependence among latent concepts, and a sparsity regularization based on the assumption that latent concepts may be sparsely activated. We implement it as follows:

$$S = ReLU(\mathbf{w}_s(\mathbf{f}_{\mathbf{x}}(\mathbf{x}) - \mathbf{b}_d) + \mathbf{b}_s), \tag{87}$$

$$\mathbf{R} = \text{ReLU}(\mathbf{w}_l(\mathbf{f}_{\mathbf{x}}(\mathbf{x}) - \mathbf{b}_d) + \mathbf{b}_l), \tag{88}$$

$$z = S + R, \tag{89}$$

$$\bar{\mathbf{f}}_{\mathbf{x}}(\mathbf{x}) = \mathbf{w}_d \mathbf{z} + \mathbf{b}_d. \tag{90}$$

Here:

- S denotes the sparse representations from the sparse encoder (with parameters  $\mathbf{w}_s$ ,  $\mathbf{b}_s$ );
- $\mathbf{R}$  denotes the structured representation from the structured encoder (with parameters  $\mathbf{w}_l, \mathbf{b}_l$ );
- z is the combined representations used for reconstruction;
- $\overline{\mathbf{f}}_{\mathbf{x}}(\mathbf{x})$  denote the reconstruction of  $\mathbf{f}_{\mathbf{x}}(\mathbf{x})$ .

The loss function for training is:

 $f = \mathbb{F} \quad \text{if } (\mathbf{y}) = \overline{\mathbf{f}} (\mathbf{y}) \|^2 + \lambda \cdot (\|\mathbf{S}\|^{p_t})$ 

$$\mathcal{L} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{train}}} \left[ \|\mathbf{f}_{\mathbf{x}}(\mathbf{x}) - \overline{\mathbf{f}}_{\mathbf{x}}(\mathbf{x})\|_{2}^{2} + \lambda_{t} \left( \|\mathbf{S}\|_{n_{t}}^{p_{t}} + \gamma \|\mathbf{R}\|_{\text{nuc}} \right) \right], \tag{91}$$

where:

- $\lambda_t$  is the dynamically adjusted sparsity coefficient at step t, following p-annealing SAE (Karvonen et al., 2024);
- $\gamma$  is a hyperparameter that balances the sparsity penalty and the low-rank regularization;
- $\|\mathbf{S}\|_{p_t}^{p_t} = \sum_i |s_i|^{p_t}$  is the adaptive  $L_{p_t}$  norm promoting sparsity, following p-annealing SAE (Karvonen et al., 2024);
- $\|\mathbf{R}\|_{nuc}$  is the nuclear norm, used to encourage low-rank structure.

We estimate the nuclear norm using the top- $k_{svd}$  singular values:

$$\|\mathbf{R}\|_{\text{nuc}} \approx \sum_{i=1}^{k_{svd}} \sigma_i,$$
 (92)

where  $\{\sigma_i\}_{i=1}^{k_{svd}}$  are the largest  $k_{svd}$  singular values, obtained via low-rank SVD (e.g., PyTorch's svd\_lowrank). The value of  $k_{svd}$  is a tunable parameter (e.g.,  $k_{svd}=64$ ) that trades off approximation accuracy and computational cost.

#### I.2 DETAILS OF EVALUATION METRIC

Obtaining  $p(c^k = 1|\mathbf{x})$  from Supervised Linear Classification. For evaluation, we first use 27 counterfactual pairs from (Park et al., 2023), also see Table 1, to train a linear classification using LogisticRegression classifier from the scikit-learn library, to obtain logits, i.e., unnormalized  $p(c^k = 1|\mathbf{x})$ , for each of the 27 pairs. As a result, for each concept in these 27 concepts, we can obtain the corresponding logit. Figure 5 shows classification accuracy of LogisticRegression classifier. Stacking these 27 logits yields the logit vector

$$\mathbf{u} = (u_1, u_2, \dots, u_{27}). \tag{93}$$

**Extracting**  $z_i$  from trained SAEs. We use the representations  $\mathbf{f_x}(\mathbf{x})$  of the same 27 counterfactual pairs from (Park et al., 2023) as input to a trained SAE, extracting the corresponding latent features  $\mathbf{z}$ . Let  $\tilde{\mathbf{z}}$  denote the element-wise exponentiation of  $\mathbf{z}$ , i.e.,  $\tilde{\mathbf{z}} = \exp(\mathbf{z})$ . This yields a feature matrix of size  $27 \times D$ , where D is the dimensionality of  $\mathbf{z}$ :

$$\tilde{\mathbf{z}} = (\tilde{\mathbf{z}}_1, \tilde{\mathbf{z}}_2, \dots, \tilde{\mathbf{z}}_{27})^T. \tag{94}$$

**Correlation Matrix and Assignment.** For the logit vector  $\mathbf{u}$ , and the feature matrix  $\tilde{\mathbf{z}}$ , we compute the Pearson correlation

$$\mathbf{R}_d = \operatorname{corr}(\mathbf{u}, \, \tilde{\mathbf{z}}_{::d}), \qquad d = 1, \dots, D \tag{95}$$

where  $\tilde{\mathbf{z}}_{:,d}$  denotes the d-th column of  $\tilde{\mathbf{z}}$ .

Note that the estimated features z from the SAE are subject to permutation indeterminacy. To address this, we apply the Hungarian algorithm to solve the assignment problem on  $\mathbf{R}_d$ . This yields the optimal assignment for each concept, allowing us to compute the assigned Pearson correlation. We report the mean Pearson correlation across the 27 concepts.

#### I.3 EXPERIMENTS AND RESULTS

We train four SAE variants—top-k SAE (Gao et al., 2025), batch-top-k SAE (Bussmann et al., 2024), p-annealing SAE (Karvonen et al., 2024), and our proposed structured SAE. Each SAE is trained three times on activations from the final hidden layer of the pretrained Pythia 70m, Pythia 1.4b, and Pythia 2.8b (Biderman et al., 2023) (download from https://huggingface.co/EleutherAI), using training data from the first 200 million tokens of the Pile corpus (Gao et al., 2020) (download from https://huggingface.co/datasets/EleutherAI/the\_pile\_deduplicated).

**Experimental setup.** We set the feature dimension of all SAE variants to be the same (D=32,768) and are trained for 20 000 optimization steps with a batch size of 10 000. We employ the Adam optimizer with an initial learning rate of  $1\times 10^{-4}$  and linearly warm up the learning rate during the first 200 steps. For the top-k and batch-top-k SAEs, we set k=32. For p-annealing SAEs, we apply a sparsity warm-up of 400 steps and an initial sparsity penalty coefficient  $\lambda_s=0.1$ . The p-annealing-LoRa SAE uses the same p-annealing settings and additionally applies a low-rank scaling factor  $\gamma=0.1$ .

**Compute resources.** All experiments are conducted on a server equipped with four NVIDIA A100 GPUs (40 GB each).

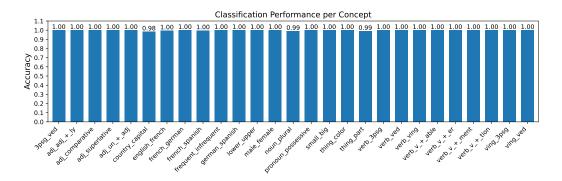


Figure 5: Classification accuracy of logistic probes across various concepts. Each bar represents the performance for a given concept.

**Pearson Correlation Coefficient.** We report Pearson Correlation Coefficient (PCC) at the 20 000-step checkpoint. As summarized in Tables 3-5, the two p-annealing variants markedly outperform the fixed-k baselines. Across the 27 concepts, the proposed structured SAE achieves better PCC across different scales of Pythia family, confirming that sparsity and the low-rank adaptation yield features that align more cleanly with human-interpretable concepts.

Table 2: Ablation results for the low-rank term scaling factor  $\gamma$ . PCC is reported as mean  $\pm$  standard deviation over multiple runs.

$\gamma$	PCC (mean $\pm$ std)
$-10^{-3}$	$0.691 \pm 0.023$
$10^{-2}$	$0.695 \pm 0.019$
$10^{-1}$	$0.704 \pm 0.013$
1	$0.685 \pm 0.021$

Ablation on the low-rank coefficient. We investigate the effect of the low-rank term scaling factor  $\gamma \in \{10^{-3}, 10^{-2}, 10^{-1}, 1\}$  while keeping all other hyperparameters fixed. As shown in Table 2, the PCC peaks at  $\gamma = 10^{-1}$ , indicating that a moderate weighting of the low-rank term is beneficial. Smaller or larger values of  $\gamma$  lead to slightly worse performance, suggesting that both under- and over-emphasis on the low-rank term can hurt the model's ability to capture the underlying structure.

Table 3: PCC results at training step 20,000 on Pythia-70m. Bold values denote the best.

PCC across different SAEs (†)				
Concepts	top-k	batch-top-k	p-annealing	Our structured SAE
3psg_ved	$0.533 \pm 0.069$	$0.466 \pm 0.012$	$0.660 \pm 0.011$	$\boldsymbol{0.684 \pm 0.008}$
adj_adj_+_ly	$0.828 \pm 0.054$	$\boldsymbol{0.918 \pm 0.005}$	$0.848 \pm 0.004$	$0.884 \pm 0.005$
adj_comparative	$0.640 \pm 0.108$	$0.646 \pm 0.087$	$0.536 \pm 0.077$	$\boldsymbol{0.695 \pm 0.070}$
adj_superlative	$0.743 \pm 0.017$	$0.757 \pm 0.028$	$0.731 \pm 0.021$	$\boldsymbol{0.792 \pm 0.033}$
adj_un_+_adj	$0.362 \pm 0.041$	$0.353 \pm 0.030$	$0.510 \pm 0.031$	$\boldsymbol{0.771 \pm 0.033}$
country_capital	$0.237 \pm 0.008$	$0.247 \pm 0.021$	$0.356 \pm 0.005$	$\boldsymbol{0.494 \pm 0.003}$
english_french	$0.501 \pm 0.076$	$0.448 \pm 0.006$	$0.629 \pm 0.008$	$\boldsymbol{0.639 \pm 0.005}$
french_german	$0.615 \pm 0.003$	$0.617 \pm 0.023$	$0.608 \pm 0.020$	$\boldsymbol{0.685 \pm 0.050}$
french_spanish	$0.627 \pm 0.069$	$0.520 \pm 0.054$	$\boldsymbol{0.643 \pm 0.033}$	$0.613 \pm 0.036$
frequent_infrequent	$0.279 \pm 0.014$	$0.292 \pm 0.017$	$0.490 \pm 0.011$	$\boldsymbol{0.509 \pm 0.018}$
german_spanish	$0.583 \pm 0.077$	$0.730 \pm 0.010$	$0.688 \pm 0.021$	$0.713 \pm 0.010$
lower_upper	$0.474 \pm 0.080$	$0.442 \pm 0.016$	$0.631 \pm 0.010$	$0.614 \pm 0.008$
male_female	$0.378 \pm 0.040$	$0.326 \pm 0.008$	$0.585 \pm 0.030$	$\boldsymbol{0.625 \pm 0.006}$
noun_plural	$0.616 \pm 0.100$	$0.501 \pm 0.038$	$\boldsymbol{0.662 \pm 0.024}$	$0.640 \pm 0.011$
pronoun_possessive	$0.897 \pm 0.035$	$0.885 \pm 0.037$	$0.967 \pm 0.022$	$\boldsymbol{0.974 \pm 0.020}$
small_big	$0.403 \pm 0.005$	$0.412 \pm 0.040$	$0.603 \pm 0.023$	$\boldsymbol{0.668 \pm 0.010}$
thing_color	$\boldsymbol{0.950 \pm 0.008}$	$0.899 \pm 0.042$	$0.918 \pm 0.020$	$0.915 \pm 0.016$
thing_part	$0.421 \pm 0.013$	$0.345 \pm 0.016$	$\boldsymbol{0.494 \pm 0.031}$	$0.506 \pm 0.008$
verb_3psg	$0.658 \pm 0.035$	$0.588 \pm 0.061$	$0.673 \pm 0.033$	$\boldsymbol{0.752 \pm 0.027}$
verb_v_+_able	$0.723 \pm 0.068$	$0.706 \pm 0.033$	$0.787 \pm 0.031$	$\boldsymbol{0.825 \pm 0.018}$
verb_v_+_er	$0.581 \pm 0.039$	$0.548 \pm 0.021$	$\boldsymbol{0.772 \pm 0.025}$	$0.748 \pm 0.011$
verb_v_+_ment	$0.750 \pm 0.013$	$\boldsymbol{0.778 \pm 0.017}$	$0.622 \pm 0.019$	$0.696 \pm 0.017$
verb_v_+_tion	$0.764 \pm 0.054$	$0.681 \pm 0.079$	$0.722 \pm 0.048$	$\boldsymbol{0.794 \pm 0.035}$
verb_ved	$0.522 \pm 0.064$	$0.583 \pm 0.087$	$0.537 \pm 0.061$	$\boldsymbol{0.668 \pm 0.052}$
verb_ving	$0.689 \pm 0.066$	$0.717 \pm 0.035$	$0.638 \pm 0.026$	$\boldsymbol{0.737 \pm 0.034}$
ving_3psg	$0.467 \pm 0.057$	$0.403 \pm 0.029$	$\boldsymbol{0.704 \pm 0.029}$	$0.704 \pm 0.010$
ving_ved	$0.557 \pm 0.043$	$0.538 \pm 0.067$	$0.691 \pm 0.053$	$\boldsymbol{0.688 \pm 0.022}$

#### J FURTHER DISCUSSION: OBSERVATIONS ON LLMS AND WORLD MODELS

# J.1 LLMs Mimic the Human World Model, Not the World Itself

As we explore the implications of our linear identifiability result, it is important to situate it within the broader context of human cognition—specifically, how humans develop and interact with an internal world model. In this subsection, we introduce the concept that LLMs Mimic the Human World Model, Not the World Itself, emphasizing the distinction between the vast physical world and the compressed abstraction humans use for reasoning and decision-making. Our analysis shows that LLMs replicate human-like abstractions through latent variable models. We further argue that LLMs aim to emulate this internal, compressed world model—rather than the physical world itself—by learning from human-generated text. This distinction provides critical insight into the success of LLMs in tasks aligned with human conceptualization and deepens our understanding of the relationship between language models and human cognition.

Humans gradually develop their understanding of the environment through learning from others and interacting with the world. This internal representation of our external environment is known as a world model. A key observation is that this model is not a direct reflection of the physical world but rather a highly compressed abstraction of it. For instance, numbers, such as 1, 2, and 3, are abstract tools created by the human mind for reasoning and problem-solving. They do not exist as tangible entities in the natural world. Similarly, many aspects of reality that escape human senses or even the most sophisticated scientific instruments are absent from our mental representations.

Interestingly, our texts reflect our mental activities and emotions, providing a window into this compressed world model. A recent study reveals a striking disparity between the limited information throughput of human behavior (approximately 10 bits/s) and the vast sensory input available (around  $10^9$  bits/s) (Zheng & Meister, 2024). This suggests that the human world model is an efficient, compressed abstraction, enabling us to reason, predict, and make decisions effectively. Despite this

Table 4: PCC results at training step 20,000 on Pythia-1.4b. Bold values denote the best.

PCC across different SAEs (†)				
Concepts	top-k	batch-top-k	p-annealing	Our structured SAE
3psg_ved	$0.577 \pm 0.015$	$0.536 \pm 0.013$	$0.629 \pm 0.012$	$\boldsymbol{0.727 \pm 0.009}$
adj_adj_+_ly	$0.870 \pm 0.013$	$0.841 \pm 0.009$	$0.897 \pm 0.020$	$0.940 \pm 0.011$
adj_comparative	$\boldsymbol{0.760 \pm 0.018}$	$0.738 \pm 0.013$	$0.640 \pm 0.011$	$0.660 \pm 0.010$
adj_superlative	$0.849 \pm 0.022$	$\boldsymbol{0.871 \pm 0.017}$	$0.767 \pm 0.013$	$0.790 \pm 0.009$
adj_un_+_adj	$0.563 \pm 0.076$	$0.556 \pm 0.107$	$0.654 \pm 0.021$	$\boldsymbol{0.749 \pm 0.012}$
average	$0.697 \pm 0.004$	$0.686 \pm 0.002$	$0.742 \pm 0.019$	$\boldsymbol{0.752 \pm 0.010}$
country_capital	$0.802 \pm 0.032$	$0.795 \pm 0.006$	$0.788 \pm 0.007$	$\boldsymbol{0.823 \pm 0.012}$
english_french	$\boldsymbol{0.836 \pm 0.014}$	$0.810 \pm 0.066$	$0.813 \pm 0.008$	$0.814 \pm 0.011$
french_german	$0.701 \pm 0.030$	$0.676 \pm 0.081$	$0.790 \pm 0.016$	$\boldsymbol{0.799 \pm 0.012}$
french_spanish	$0.737 \pm 0.006$	$0.706 \pm 0.041$	$\boldsymbol{0.823 \pm 0.010}$	$0.762 \pm 0.013$
frequent_infrequent	$0.522 \pm 0.029$	$0.569 \pm 0.003$	$0.588 \pm 0.012$	$\boldsymbol{0.646 \pm 0.011}$
german_spanish	$0.687 \pm 0.023$	$0.673 \pm 0.038$	$\boldsymbol{0.796 \pm 0.009}$	$0.783 \pm 0.014$
lower_upper	$\boldsymbol{0.765 \pm 0.007}$	$0.746 \pm 0.007$	$0.716 \pm 0.010$	$0.726 \pm 0.012$
male_female	$0.584 \pm 0.047$	$\boldsymbol{0.617 \pm 0.011}$	$0.593 \pm 0.006$	$0.546 \pm 0.007$
noun_plural	$0.688 \pm 0.006$	$0.703 \pm 0.029$	$\boldsymbol{0.851 \pm 0.014}$	$0.814 \pm 0.010$
pronoun_possessive	$0.937 \pm 0.018$	$0.869 \pm 0.031$	$0.980 \pm 0.012$	$\boldsymbol{0.989 \pm 0.008}$
small_big	$0.287 \pm 0.007$	$0.300 \pm 0.008$	$0.536 \pm 0.013$	$0.533 \pm 0.010$
thing_color	$0.913 \pm 0.019$	$\boldsymbol{0.937 \pm 0.009}$	$0.928 \pm 0.012$	$0.904 \pm 0.011$
thing_part	$0.397 \pm 0.005$	$0.402 \pm 0.013$	$0.440 \pm 0.010$	$\boldsymbol{0.531 \pm 0.014}$
verb_3psg	$0.675 \pm 0.090$	$0.619 \pm 0.011$	$0.709 \pm 0.008$	$\boldsymbol{0.716 \pm 0.012}$
verb_v_+_able	$0.785 \pm 0.048$	$0.743 \pm 0.026$	$\boldsymbol{0.899 \pm 0.015}$	$0.798 \pm 0.013$
verb_v_+_er	$0.707 \pm 0.012$	$0.745 \pm 0.019$	$0.752 \pm 0.011$	$\boldsymbol{0.781 \pm 0.012}$
verb_v_+_ment	$0.769 \pm 0.014$	$0.603 \pm 0.033$	$0.831 \pm 0.010$	$0.741 \pm 0.009$
verb_v_+_tion	$0.857 \pm 0.025$	$0.879 \pm 0.013$	$\boldsymbol{0.841 \pm 0.012}$	$0.818 \pm 0.010$
verb_ved	$0.599 \pm 0.027$	$0.674 \pm 0.032$	$0.682 \pm 0.009$	$\boldsymbol{0.710 \pm 0.008}$
verb_ving	$0.730 \pm 0.004$	$0.721 \pm 0.053$	$\boldsymbol{0.800 \pm 0.013}$	$0.794 \pm 0.011$
ving_3psg	$0.612 \pm 0.064$	$0.625 \pm 0.036$	$0.642 \pm 0.010$	$\boldsymbol{0.745 \pm 0.014}$
ving_ved	$0.605 \pm 0.019$	$0.558\pm0.008$	$0.645 \pm 0.011$	$\boldsymbol{0.660 \pm 0.012}$

compression, humans have flourished as the dominant species on Earth, demonstrating the power of such a streamlined model.

LLMs aim to mimic not the vast and unbounded world but this human-compressed world model, which is significantly smaller and more manageable. By learning from human text, which encodes this abstraction, LLMs effectively replicate the patterns, reasoning, and abstractions that have proven successful for humans. This explains the impressive performance of LLMs in tasks that align with human understanding. Furthermore, the overlap between the latent space of LLMs (representing human concepts) and the observed space (human text) provides a powerful mechanism for aligning human-like abstractions with model predictions. For instance, modifying words in the observed space often corresponds to predictable changes in latent concepts.

#### J.2 LLMs Versus Pure Vision Models: A Fundamental Difference

While LLMs model human language, which has already undergone significant compression through human cognition, vision models face a fundamentally different challenge. In the previous subsection, we discussed how LLMs mimic the human world model, leveraging compressed abstractions derived from human-generated text. In contrast, vision models operate on raw, high-dimensional data from visual inputs. Moreover, the training data for vision models represents only a tiny fraction of the universe, constrained by the limitations of capturing devices and datasets. This vastness and lack of compression make the task of building generalizable representations in vision models inherently more complex than in NLP-based LLMs.

This difference in the nature of their training data and observed space might explain the behavioral differences between LLMs and pure vision models. While LLMs benefit from the inherent abstraction and compression of human language, vision models must contend with raw, unprocessed inputs that

Table 5: PCC results at training step 20,000 on Pythia-2.8b. Bold values denote the best.

PCC across different SAEs (†)				
Concepts	top-k	batch-top-k	p-annealing	Our structured SAE
3psg_ved	$0.589 \pm 0.025$	$0.513 \pm 0.024$	$0.645 \pm 0.012$	$\boldsymbol{0.734 \pm 0.027}$
adj_adj_+_ly	$0.869 \pm 0.036$	$0.797 \pm 0.050$	$0.919 \pm 0.020$	$0.941\pm0.013$
adj_comparative	$0.847 \pm 0.020$	$\boldsymbol{0.859 \pm 0.010}$	$0.817 \pm 0.015$	$0.757 \pm 0.028$
adj_superlative	$0.885 \pm 0.011$	$0.897 \pm 0.010$	$\boldsymbol{0.901 \pm 0.017}$	$0.787 \pm 0.009$
adj_un_+_adj	$0.631 \pm 0.086$	$0.690 \pm 0.033$	$0.590 \pm 0.012$	$\boldsymbol{0.693 \pm 0.021}$
average	$0.707 \pm 0.008$	$0.704 \pm 0.002$	$0.762 \pm 0.011$	$\boldsymbol{0.771 \pm 0.022}$
country_capital	$0.810 \pm 0.040$	$0.791 \pm 0.034$	$\boldsymbol{0.892 \pm 0.014}$	$0.832 \pm 0.020$
english_french	$\boldsymbol{0.854 \pm 0.010}$	$0.834 \pm 0.010$	$0.842 \pm 0.015$	$0.847 \pm 0.027$
french_german	$0.659 \pm 0.060$	$0.710 \pm 0.035$	$0.743 \pm 0.016$	$\boldsymbol{0.799 \pm 0.022}$
french_spanish	$0.741 \pm 0.027$	$0.711 \pm 0.012$	$\boldsymbol{0.828 \pm 0.011}$	$0.794 \pm 0.019$
frequent_infrequent	$0.523 \pm 0.066$	$0.576 \pm 0.036$	$0.572 \pm 0.018$	$\boldsymbol{0.670 \pm 0.028}$
german_spanish	$0.643 \pm 0.040$	$0.646 \pm 0.063$	$\boldsymbol{0.816 \pm 0.021}$	$0.805 \pm 0.030$
lower_upper	$\boldsymbol{0.772 \pm 0.027}$	$0.758 \pm 0.010$	$0.704 \pm 0.012$	$0.739 \pm 0.025$
male_female	$0.640 \pm 0.040$	$0.570 \pm 0.014$	$\boldsymbol{0.721 \pm 0.018}$	$0.640 \pm 0.020$
noun_plural	$0.761 \pm 0.024$	$0.737 \pm 0.014$	$\boldsymbol{0.879 \pm 0.023}$	$0.830 \pm 0.027$
pronoun_possessive	$0.918 \pm 0.049$	$0.914 \pm 0.011$	$0.980 \pm 0.015$	$\boldsymbol{0.995 \pm 0.021}$
small_big	$0.347 \pm 0.007$	$0.400 \pm 0.043$	$0.481 \pm 0.014$	$\boldsymbol{0.602 \pm 0.032}$
thing_color	$0.918 \pm 0.019$	$\boldsymbol{0.932 \pm 0.012}$	$0.923 \pm 0.007$	$0.901 \pm 0.015$
thing_part	$0.373 \pm 0.033$	$0.319 \pm 0.021$	$0.457 \pm 0.010$	$\boldsymbol{0.506 \pm 0.025}$
verb_3psg	$0.523 \pm 0.071$	$0.578 \pm 0.025$	$0.643 \pm 0.012$	$\boldsymbol{0.739 \pm 0.021}$
verb_v_+_able	$0.751 \pm 0.022$	$0.736 \pm 0.023$	$\boldsymbol{0.832 \pm 0.018}$	$0.828 \pm 0.026$
verb_v_+_er	$0.715 \pm 0.052$	$0.790 \pm 0.011$	$\boldsymbol{0.832 \pm 0.025}$	$0.806 \pm 0.022$
verb_v_+_ment	$0.784 \pm 0.036$	$0.762 \pm 0.032$	$\boldsymbol{0.860 \pm 0.015}$	$0.854 \pm 0.028$
verb_v_+_tion	$0.805 \pm 0.053$	$0.845 \pm 0.085$	$\boldsymbol{0.864 \pm 0.023}$	$0.861 \pm 0.030$
verb_ved	$0.720 \pm 0.041$	$0.728 \pm 0.039$	$0.664 \pm 0.012$	$\boldsymbol{0.752 \pm 0.018}$
verb_ving	$0.732 \pm 0.060$	$0.755 \pm 0.035$	$\boldsymbol{0.845 \pm 0.020}$	$0.819 \pm 0.027$
ving_3psg	$\boldsymbol{0.728 \pm 0.082}$	$0.559 \pm 0.117$	$0.690 \pm 0.015$	$0.675 \pm 0.022$
ving_ved	$0.562 \pm 0.022$	$0.589 \pm 0.009$	$\boldsymbol{0.624 \pm 0.017}$	$0.620 \pm 0.019$

require far greater generalization capabilities. This underscores the importance of understanding the nuances of each modality when designing and evaluating AI systems.

# K More Results on Llama-3 and DeepSeek-R1

We conduct additional experiments on recent LLMs, including Llama-3 and DeepSeek-R1, to further evaluate our findings. The experimental setup strictly adheres to the settings described in Section H.2. This enables a comprehensive investigation of our findings, ensuring that the results are thoroughly evaluated and validated across diverse LLM architectures and experimental conditions. Overall, we can see, the product  $\mathbf{A}_s \times \mathbf{W}_s$  approximates the identity matrix, supporting the theoretical findings outlined in Corollary G.2.

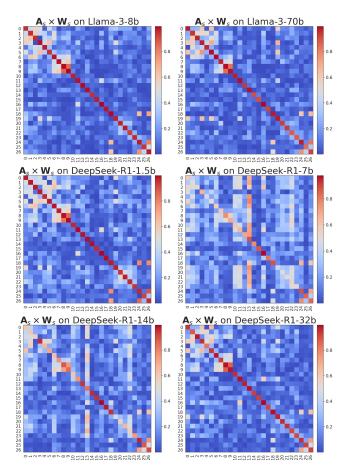


Figure 6: Results of the product  $\mathbf{A}_s \times \mathbf{W}_s$  across the LLaMA-3 and DeepSeek-R1 model families. Here,  $\mathbf{A}_s$  represents a matrix derived from the feature differences of 27 counterfactual pairs, while  $\mathbf{W}_s$  is a weight matrix obtained from a linear classifier trained on these features.

# L FUTURE DIRECTIONS

 Rethinking Invertibility Assumptions in Causal Representation Learning Our identifiability analysis is closely related to the concept of identifiability in causal representation learning. Notably, to the best of our knowledge, this is the first work to explore approximate identifiability in the context of *non-invertible* mappings from latent space to observed space—a departure from the commonly upheld *invertibility* assumption in the causal representation learning community (Brehmer et al., 2022; Von Kügelgen et al., 2021; Massidda et al., 2023; von Kügelgen et al., 2023; Ahuja et al., 2023; Seigal et al., 2022; Shen et al., 2022; Liu et al., 2022). We hope that our work will inspire future research aimed at overcoming the limitations imposed by invertibility assumptions in causal representation learning.

Embedding Causal Reasoning in LLMs Through Linear Unmixing Our linear identifiability result lays a foundation for uncovering latent causal relationships among concepts, especially when these variables exhibit causal dependencies within the proposed latent variable model. By showing that the representations in LLMs are linear mixtures of latent causal variables, our analysis shows that linear unmixing of these representations may allow for the identification of underlying latent causal structures. This approach not only opens up the possibility of understanding causal dynamics within LLMs but also suggests that causal reasoning, particularly in latent spaces, could be achievable through the exploration of these linear unmixing techniques. We believe this work marks a pivotal step toward embedding robust causal reasoning capabilities into LLMs.

# M ACKNOWLEDGMENT OF LLMS USAGE

We disclose that large language models (LLMs) were used only to correct typos, improve grammar, and refine phrasing. All scientific contributions, including problem formulation, methodology design, theoretical proofs, experiments, and analysis of results, are exclusively the work of the authors.