Beyond the Black Box: Identifiable Interpretation and Control in Generative Models via Causal Minimality

Lingjing Kong*1, Shaoan Xie*1, Guangyi Chen 1,2 , Yuewen Sun 1,2 , Xiangchen Song 1 , Eric P. Xing 1,2 , Kun Zhang 1,2

*Equal contribution

¹Carnegie Mellon University, Pittsburgh, PA, USA

²Mohamed bin Zayed University of Artificial Intelligence, Abu Dhabi, UAE

Abstract

Deep generative models, while revolutionizing fields like image and text generation, largely operate as opaque "black boxes", hindering human understanding, control, and alignment. While methods like sparse autoencoders (SAEs) show remarkable empirical success, they often lack theoretical guarantees, risking subjective insights. Our primary objective is to establish a principled foundation for interpretable generative models. We demonstrate that the principle of causal minimality – favoring the simplest causal explanation – can endow the latent representations of diffusion vision and autoregressive language models with clear causal interpretation and robust, component-wise identifiable control. We introduce a novel theoretical framework for hierarchical selection models, where higher-level concepts emerge from the constrained composition of lower-level variables, better capturing the complex dependencies in data generation. Under theoretically derived minimality conditions (manifesting as sparsity or compression constraints), we show that learned representations can be equivalent to the true latent variables of the data-generating process. Empirically, applying these constraints to leading generative models allows us to extract their innate hierarchical concept graphs, offering fresh insights into their internal knowledge organization. Furthermore, these causally grounded concepts serve as levers for fine-grained model steering, paving the way for transparent, reliable systems.

1 Introduction

Deep generative models, such as diffusion [25, 56] and language models [4], are reshaping numerous domains. However, their complexity often renders them 'black boxes, hindering understanding and control. While empirical tools like sparse autoencoders (SAEs) offer methods for probing these models, they lack theoretical guarantees, risking subjective interpretations. In this work, we tackle this challenge, seeking a principled foundation for interpretable and controllable generative models.

We identify *causal minimality* [55, 65, 24] as the formal principle connecting practices like sparsity to the recovery of meaningful, hierarchical concepts. We apply this to text-to-image (T2I) diffusion models and autoregressive LMs, finding that sparsity constraints are instrumental for identifying visual and textual concepts.

A cornerstone of our contribution is establishing the first identifiability results for *selection*-based [81, 66, 23, 76, 2, 14, 9, 5] hierarchical models, where higher-level concepts emerge as effects of lower-level compositions. This diverges from traditional models where influence propagates from high to low levels [54, 7, 78] and, we argue, better captures the dependencies in concept formation (e.g., a "car" from its parts). Our framework establishes *component-wise* identifiability for these

39th Conference on Neural Information Processing Systems (NeurIPS 2025) Workshop: Mechanistic Interpretability.

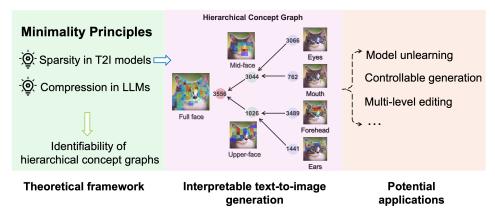


Figure 1: Our causal minimality principle enables interpretable text-to-image generation through hierarchical concept graphs, with implications for downstream tasks.

models (Conditions 3.2-iv and A.1-iii), proving that under these minimality conditions, learned representations can be equivalent to the true latent variables. Empirically, this allows us to extract and steer innate hierarchical concept graphs (Figure 1), offering new levers for model control. Due to page limits, we focus on T2I models in the main text and defer LMs to Appendix A.

2 Deep Generative Models as Hierarchical Concept Models

Notations. We denote random variables with upper-case characters (e.g., X) and values with lower-case characters (e.g., x). We use bold fonts for multidimensional objects (e.g., x) with dimensionality $n(\cdot)$. Parents $\mathrm{Pa}(\cdot)$ and children $\mathrm{Ch}(\cdot)$ relations are defined based on the selection graph (Figure 2). We denote $[M] := \{1,\ldots,M\}$. We denote the image as $\mathbf{X} \in \mathbb{R}^{n(\mathbf{X})}$, text as $\mathbf{D} \in \mathbb{N}^{n(\mathbf{D})}$, and visual concepts as $\mathbf{Z} := [\mathbf{Z}_1, \cdots, \mathbf{Z}_{L_V}]$, where $\mathbf{Z}_l \in \mathbb{R}^{n(\mathbf{Z}_l)}$ are concepts at level l.

Hierarchical processes and selection mechanisms. Our framework conceptualizes high-level concepts

Our framework conceptualizes high-level concepts as emerging from lower-level ones via a *selection mechanism* [81, 66, 23, 76, 2, 14, 9, 5]. A higher-level concept V_l is an effect of its constituents V_{l+1} (i.e., its "parents"): $V_l := g_{V_l}(V_{l+1})$. This "selection" perspective is critical for modeling how abstract concepts (e.g., "bicycle") enforce coherence among their parts (e.g., "wheels", "frame"). Traditional models [7, 53, 78] often require dense intralevel graphs to capture these dependencies, violating the minimality principle. The generative pathway inverts this process, sampling from abstract concepts to concrete details:

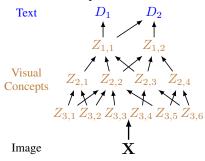


Figure 2: A visual concept graph. We denote text as **D**, visual concepts as **Z**, and the image as **X**. See Appendix A for the text counterpart.

$$\mathbf{Z}_0 \sim \mathbb{P}(\mathbf{Z}_0), \quad \mathbf{Z}_l \sim \mathbb{P}(\mathbf{Z}_l | \mathbf{Z}_{l-1}), \quad l \in \{1, \dots, L_{V} + 1\},$$
 (1)

where we denote $\mathbf{Z}_0 := \mathbf{D}$ and $\mathbf{Z}_{L_V+1} := \mathbf{X}$.

Connections to text-to-image diffusion models. This hierarchical generation aligns with the iterative denoising in diffusion models [63, 25]. Each step f_t can be seen as an autoencoder extracting a representation $\mathbf{Z}_{\mathcal{S}(t)}$. Representations from high-noise levels (large t) correspond to abstract concepts (e.g., \mathbf{Z}_l with small l), while low-noise levels capture concrete details (large l). The text prompt \mathbf{D} guides the most abstract concepts (e.g., \mathbf{Z}_l). This motivates our core question: are these learned representations *identifiable* with the true latent concepts?

3 Identifiable Representations under Causal Minimality

We posit that by adhering to the *causal minimality principle* [55, 65, 24], we can achieve component-wise identifiability. For visual concepts, minimality manifests as *sparse graphical dependencies*. For discrete text concepts, it translates to seeking the *most compressed representation* (see Appendix A). We focus on *hierarchical selection models*, where high-level concepts are effects of

lower-level ones (colliders), rendering many existing identifiability results inapplicable. Our goal is *component-wise identifiability*:

Definition 3.1 (Component-wise Identifiability). Let \mathbf{Z} and $\hat{\mathbf{Z}}$ be variables under two model specifications. We say that \mathbf{Z} and $\hat{\mathbf{Z}}$ are *identified component-wise* if there exists a permutation π such that for each $i \in [n(\mathbf{Z})]$, $\hat{Z}_i = h_i(Z_{\pi(i)})$ where h_i is an invertible function.

This strong form ensures a one-to-one mapping between learned and true concepts, which is vital for unambiguous interpretation. We assume the standard faithfulness condition [64]. In the following, we consider the identification of continuous latent visual concepts **Z**.

Condition 3.2 (Visual Concept Identification Conditions).

- i Informativeness: There exists a diffeomorphism $g_l: (\mathbf{Z}_l, \epsilon_l) \mapsto \mathbf{X}$ for $l \in [0, L]$, where ϵ_l denotes independent exogenous variables.
- ii Smooth Density: The probability density function $p(\mathbf{z}_{l+1}|\mathbf{z}_l)$ is smooth for any $l \in [L_V]$.
- iii Sufficient Variability: For each Z and its parents $\tilde{\mathbf{Z}} := \operatorname{Pa}(Z)$, at any value $\tilde{\mathbf{z}}$ of $\tilde{\mathbf{Z}}$, there exist $n(\tilde{\mathbf{Z}}) + 1$ distinct values of Z, denoted as $\{z^{(n)}\}_{n=0}^{n(\tilde{Z})}$, such that the vectors $\mathbf{w}(\tilde{\mathbf{z}}, z^{(n)}) \mathbf{w}(\tilde{\mathbf{z}}, z^{(0)})$ are linearly independent where $\mathbf{w}(\tilde{\mathbf{z}}, z) = \left(\frac{\partial \log p(\tilde{\mathbf{z}}|z)}{\partial \tilde{z}_1}, \dots, \frac{\partial \log p(\tilde{\mathbf{z}}|z)}{\partial \tilde{z}_{n(\tilde{\mathbf{z}})}}\right)$.
- iv Sparse Connectivity (Minimality): For each parent concept \tilde{Z} , there exists a subset of its children $\mathbf{Z} \subseteq \operatorname{Ch}(\tilde{Z})$ such that their only common parent is \tilde{Z} , i.e., $\bigcap_{Z \in \mathbf{Z}} \operatorname{Pa}(Z) = \{\tilde{Z}\}$.

Interpreting Condition 3.2. Conds. 3.2-i (Informativeness) and 3.2-ii (Smooth Density) are standard regularity assumptions [29, 30, 33, 32, 70, 39]. 3.2-iii (Sufficient Variability) ensures concepts respond distinctly. 3.2-iv (Sparse Connectivity) is our causal minimality constraint, positing a sparse causal graph with unique connectivity "fingerprints," crucial for disentanglement [80, 44, 73, 45, 43].

Theorem 3.3 (Visual Concept Identification). Assume the process for visual concepts in (1). If a model specification θ_V satisfies Condition 3.2, and an alternative specification $\hat{\theta}_V$ satisfies Conditions 3.2-i and 3.2-ii, along with a sparsity constraint such that for corresponding \hat{Z} and Z:

$$n(\operatorname{Pa}(\hat{Z})) \le n(\operatorname{Pa}(Z)),$$
 (2)

then, if both models θ_V and $\hat{\theta}_V$ generate the same observed data distribution $\mathbb{P}(\mathbf{X})$, the latent visual concepts \mathbf{Z}_l are component-wise identifiable for every level $l \in [L_V]$.

Proof sketch for Theorem 3.3. The proof (full details in Appendix C) proceeds by identifying the hierarchy level-by-level, from top to bottom. 1) Paired text data \mathbf{D} acts as an auxiliary variable for \mathbf{Z}_1 , allowing identification of influenced subspaces [29, 30, 41]. 2) Given Condition 3.2-iv, the intersections of these subspaces can be identified component-wise [70, 74, 40]. 3) Once \mathbf{Z}_1 is identified, it serves as the auxiliary variable to identify \mathbf{Z}_2 , and this process is repeated.

Implications for text-to-image diffusion models. Theorem 3.3 underscores that the *sparsity constraint* (2) *is pivotal for identifying true visual concepts*. In practice, this constraint can be encouraged by techniques like SAEs applied to the internal representations (e.g., U-Net features at various timesteps) of diffusion models. By sparsifying activations, SAEs promote the sparser conceptual graphs required by our theory. Our experiments in Section 4 demonstrate this.

4 Experiments

We demonstrate our method's hierarchical interpretation on T2I models. Full empirical results, implementation details, and ablations are in Appendix F.

Hierarchical concept graph visualization. Figure 3 shows a hierarchical concept graph learned by our method. Nodes from early timesteps (e.g., 899, brown) represent high-level concepts (node 3556, "full cat face"). Nodes from intermediate timesteps (e.g., 500, green) capture mid-level regions (node 3044, "central face"), while nodes from late timesteps (e.g., 100, blue) encode fine details (node 3066, "eyes"). This clearly illustrates a coarse-to-fine hierarchy.

Concept steering in hierarchical graphs. We perform concept steering by modifying the SAE feature $x' = x + \lambda D(E(x))$ and feeding it back into the diffusion process. As shown in Fig. 3

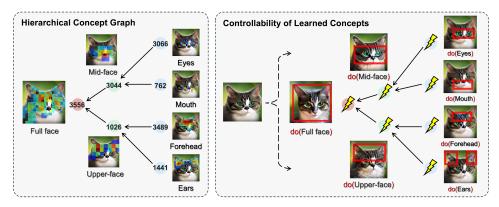


Figure 3: **Hierarchical concept graph for a text-to-image model.** Our method recovers meaningful hierarchical structures. On the right, we demonstrate feature steering: intervening on a high-level concept ("Full face") alters the cat's entire facial structure, while intervening on a lower-level concept (e.g., "Eye") produces a localized edit. More examples in Appendix F.

Method	I2P↓	R	ING-A	-BELI	,	P4D↓	UATK ↓	CC	OCO
		K77	K38	K16	AVG			FID↓	CLIP↑
SD 1.4	17.8	85.26	87.37	93.68	88.10	98.70	69.70	16.71	31.3
ESD	2.87	20.00	29.47	35.79	28.42	15.49	2.87	18.18	30.2
SA	2.81	63.15	56.84	56.84	58.94	12.68	2.81	25.80	29.7
CA	1.04	86.32	91.69	94.26	90.76	5.63	1.04	24.12	30.1
MACE	1.51	2.10	0.00	0.00	0.70	2.82	1.51	16.80	28.7
UCE	0.87	10.52	9.47	12.61	10.87	9.86	0.87	17.99	30.2
RECE	0.72	5.26	4.21	5.26	4.91	5.63	0.72	17.74	30.2
SDID	3.77	94.74	95.79	90.53	93.68	69.54	30.99	22.16	31.1
SLD-MAX	1.74	23.16	32.63	42.11	32.63	9.14	2.44	28.75	28.4
SD1.4-NegPrompt	0.74	17.89	40.42	34.74	31.68	10.00	1.46	18.33	30.1
SAFREE	1.45	35.78	47.36	55.78	46.31	10.56	1.45	19.32	30.1
TRASCE	0.45	1.05	2.10	2.10	1.75	3.97	0.70	17.41	29.9
ConceptSteer	0.36	3.16	8.42	9.47	7.02	1.99	2.11	18.67	30.8
Ours	0.25	1.05	0.00	2.11	1.05	0.66	2.11	17.02	31.3

Table 1: **Model unlearning comparisons**. Our method delivers competitive results on unlearning tasks without compromising standard text-to-image generation.

(right), steering high-level node 3556 alters the entire cat face, while steering node 1026 adjusts only the upper region, demonstrating localized, hierarchical control.

Model unlearning. Table 1 reports quantitative results on four nudity-removal benchmarks [58, 69, 6, 79]. Our method achieves state-of-the-art performance. To verify generalization, we apply feature steering on MSCOCO [47] prompts; the resulting low FID and high CLIP scores indicate effective unlearning without degrading overall text-to-image quality.

Ablation. As established in our theory, sparsity is crucial for identifiability. Fig. 6 empirically validates this. When sparsity is not enforced (left), the resulting graph is overly dense and uninterpretable. Conversely, imposing excessive sparsity (right) leads to an overly pruned graph that fails to capture the generation process. This highlights the importance of a balanced sparsity, as predicted by our theory, to achieve an interpretable and meaningful concept hierarchy.

5 Conclusion

In this work, we present a theoretical framework using causal minimality for identifying latent concepts in hierarchical selection models. We prove that generative model representations can map to true latent variables. Empirically, applying these constraints enables extracting meaningful hierarchical concept graphs from leading models, enhancing interpretability and grounded control.

References

- [1] Animashree Anandkumar, Daniel Hsu, Adel Javanmard, and Sham Kakade. Learning linear bayesian networks with latent variables. In *International Conference on Machine Learning*, pages 249–257. PMLR, 2013.
- [2] Elias Bareinboim, Jin Tian, and Judea Pearl. Recovering from selection bias in causal and statistical inference. *Probabilistic and causal inference: The works of Judea Pearl*, pages 433–450, 2022.
- [3] Joseph Bloom, Curt Tigges, Anthony Duong, and David Chanin. Saelens. https://github.com/jbloomAus/SAELens, 2024.
- [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- [5] Leihao Chen, Onno Zoeter, and Joris M Mooij. Modeling latent selection with structural causal models. *arXiv preprint arXiv:2401.06925*, 2024.
- [6] Zhi-Yi Chin, Chieh-Ming Jiang, Ching-Chun Huang, Pin-Yu Chen, and Wei-Chen Chiu. Prompting4debugging: Red-teaming text-to-image diffusion models by finding problematic prompts. *arXiv preprint arXiv:2309.06135*, 2023.
- [7] Myung Jin Choi, Vincent YF Tan, Animashree Anandkumar, and Alan S Willsky. Learning latent tree graphical models. *Journal of Machine Learning Research*, 12:1771–1812, 2011.
- [8] Joel E Cohen and Uriel G Rothblum. Nonnegative ranks, decompositions, and factorizations of nonnegative matrices. *Linear Algebra and its Applications*, 190:149–168, 1993.
- [9] Juan D Correa, Jin Tian, and Elias Bareinboim. Identification of causal effects in the presence of selection bias. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 2744–2751, 2019.
- [10] Thomas M Cover. Elements of information theory. John Wiley & Sons, 1999.
- [11] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023.
- [12] Bartosz Cywiński and Kamil Deja. Saeuron: Interpretable concept unlearning in diffusion models with sparse autoencoders. *arXiv preprint arXiv:2501.18052*, 2025.
- [13] Xinshuai Dong, Biwei Huang, Ignavier Ng, Xiangchen Song, Yujia Zheng, Songyao Jin, Roberto Legaspi, Peter Spirtes, and Kun Zhang. A versatile causal discovery framework to allow causally-related hidden variables. In *The Twelfth International Conference on Learning Representations*, 2023.
- [14] Patrick Forré and Joris M Mooij. Causal calculus in the presence of cycles, latent confounders and selection bias. In *Uncertainty in Artificial Intelligence*, pages 71–80. PMLR, 2020.
- [15] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2426–2436, 2023.
- [16] Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5111–5120, 2024.

- [17] 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.
- [18] Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. arXiv preprint arXiv:2406.04093, 2024.
- [19] Chao Gong, Kai Chen, Zhipeng Wei, Jingjing Chen, and Yu-Gang Jiang. Reliable and efficient concept erasure of text-to-image diffusion models. In *European Conference on Computer Vision*, pages 73–88. Springer, 2024.
- [20] Yuqi Gu and David B. Dunson. Bayesian pyramids: Identifiable multilayer discrete latent structure models for discrete data. In *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 2023.
- [21] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. Advances in neural information processing systems, 33:9841–9850, 2020.
- [22] Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep generative models. Advances in Neural Information Processing Systems, 36:17170– 17194, 2023.
- [23] Miguel A Hernán, Sonia Hernández-Díaz, and James M Robins. A structural approach to selection bias. *Epidemiology*, 15(5):615–625, 2004.
- [24] Christopher Hitchcock. Probabilistic Causation. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Spring 2021 edition, 2021.
- [25] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems 33 (NeurIPS 2020)*, 2020.
- [26] Biwei Huang, Charles Jia Han Low, Feng Xie, Clark Glymour, and Kun Zhang. Latent hierarchical causal structure discovery with rank constraints. Advances in Neural Information Processing Systems, 35:5549–5561, 2022.
- [27] Victor Shea-Jay Huang, Le Zhuo, Yi Xin, Zhaokai Wang, Peng Gao, and Hongsheng Li. Tide: Temporal-aware sparse autoencoders for interpretable diffusion transformers in image generation. *arXiv preprint arXiv:2503.07050*, 2025.
- [28] 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*, 2023.
- [29] Aapo Hyvarinen and Hiroshi Morioka. Unsupervised feature extraction by time-contrastive learning and nonlinear ica. *Advances in neural information processing systems*, 29, 2016.
- [30] 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*, pages 859–868. PMLR, 2019.
- [31] Anubhav Jain, Yuya Kobayashi, Takashi Shibuya, Yuhta Takida, Nasir Memon, Julian Togelius, and Yuki Mitsufuji. Trasce: Trajectory steering for concept erasure. arXiv preprint arXiv:2412.07658, 2024.
- [32] Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoencoders and nonlinear ica: A unifying framework. In *International Conference on Artificial Intelligence and Statistics*, pages 2207–2217. PMLR, 2020.
- [33] Ilyes Khemakhem, Ricardo Monti, Diederik Kingma, and Aapo Hyvarinen. Ice-beem: Identifiable conditional energy-based deep models based on nonlinear ica. *Advances in Neural Information Processing Systems*, 33:12768–12778, 2020.

- [34] Dahye Kim and Deepti Ghadiyaram. Concept steerers: Leveraging k-sparse autoencoders for controllable generations. *arXiv preprint arXiv:2501.19066*, 2025.
- [35] Dahye Kim, Xavier Thomas, and Deepti Ghadiyaram. Revelio: Interpreting and leveraging semantic information in diffusion models. *arXiv preprint arXiv:2411.16725*, 2024.
- [36] 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.
- [37] Bohdan Kivva, Goutham Rajendran, Pradeep Ravikumar, and Bryon Aragam. Identifiability of deep generative models without auxiliary information. *Advances in Neural Information Processing Systems*, 35:15687–15701, 2022.
- [38] Lingjing Kong, Guangyi Chen, Biwei Huang, Eric P. Xing, Yuejie Chi, and Kun Zhang. Learning discrete concepts in latent hierarchical models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [39] Lingjing Kong, Biwei Huang, Feng Xie, Eric Xing, Yuejie Chi, and Kun Zhang. Identification of nonlinear latent hierarchical models. Advances in Neural Information Processing Systems, 36, 2023.
- [40] Lingjing Kong, Martin Q. Ma, Guangyi Chen, Eric P. Xing, Yuejie Chi, Louis-Philippe Morency, and Kun Zhang. Understanding masked autoencoders via hierarchical latent variable models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7918–7928, June 2023.
- [41] Lingjing Kong, Shaoan Xie, Weiran Yao, Yujia Zheng, Guangyi Chen, Petar Stojanov, Victor Akinwande, and Kun Zhang. Partial disentanglement for domain adaptation. In *International Conference on Machine Learning*, pages 11455–11472. PMLR, 2022.
- [42] Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu. Ablating concepts in text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22691–22702, 2023.
- [43] Sébastien Lachapelle, Tristan Deleu, Divyat Mahajan, Ioannis Mitliagkas, Yoshua Bengio, Simon Lacoste-Julien, and Quentin Bertrand. Synergies between disentanglement and sparsity: Generalization and identifiability in multi-task learning. arXiv preprint arXiv:2211.14666, 2022.
- [44] Sébastien Lachapelle, Pau Rodríguez López, Yash Sharma, Katie Everett, Rémi Le Priol, Alexandre Lacoste, and Simon Lacoste-Julien. Nonparametric partial disentanglement via mechanism sparsity: Sparse actions, interventions and sparse temporal dependencies. *arXiv* preprint arXiv:2401.04890, 2024.
- [45] Sébastien Lachapelle, Pau Rodriguez, Yash Sharma, Katie E Everett, Rémi Le Priol, Alexandre Lacoste, and Simon Lacoste-Julien. Disentanglement via mechanism sparsity regularization: A new principle for nonlinear ica. In *Conference on Causal Learning and Reasoning*, pages 428–484. PMLR, 2022.
- [46] Hang Li, Chengzhi Shen, Philip Torr, Volker Tresp, and Jindong Gu. Self-discovering interpretable diffusion latent directions for responsible text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12006–12016, 2024.
- [47] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [48] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv* preprint *arXiv*:2211.01095, 2022.

- [49] Shilin Lu, Zilan Wang, Leyang Li, Yanzhu Liu, and Adams Wai-Kin Kong. Mace: Mass concept erasure in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6430–6440, 2024.
- [50] Gemma E Moran, Dhanya Sridhar, Yixin Wang, and David M Blei. Identifiable variational autoencoders via sparse decoding. *arXiv preprint arXiv:2110.10804*, 2021.
- [51] Anish Mudide, Joshua Engels, Eric J Michaud, Max Tegmark, and Christian Schroeder de Witt. Efficient dictionary learning with switch sparse autoencoders. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [52] 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.
- [53] J. Pearl. Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference. Morgan Kaufmann, 1988.
- [54] Judea Pearl. Causality. Cambridge university press, 2009.
- [55] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Elements of causal inference: foundations and learning algorithms. The MIT Press, 2017.
- [56] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR 2022)*, 2022.
- [57] 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.
- [58] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22522–22531, 2023.
- [59] Christoph Schuhmann, Andreas Köpf, Theo Coombes, Richard Vencu, Romain Beaumont, and Benjamin Trom. Laion-coco: 600m synthetic captions from laion2b-en. LAION.ai blog, September 2022.
- [60] Claude E Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948.
- [61] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. Interfacegan: Interpreting the disentangled face representation learned by gans. *IEEE transactions on pattern analysis and machine intelligence*, 44(4):2004–2018, 2020.
- [62] Wei Shi, Sihang Li, Tao Liang, Mingyang Wan, Gojun Ma, Xiang Wang, and Xiangnan He. Route sparse autoencoder to interpret large language models. *CoRR*, abs/2503.08200, March 2025.
- [63] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the International Conference on Machine Learning (ICML 2015)*, 2015.
- [64] Peter Spirtes, Clark Glymour, and Richard Scheines. Causation, prediction, and search. MIT press, 2001.
- [65] Peter Spirtes, Clark N Glymour, and Richard Scheines. Causation, Prediction, and Search. MIT press, 2000.
- [66] Peter Spirtes, Christopher Meek, and Thomas Richardson. Causal inference in the presence of latent variables and selection bias. In *Proceedings of the Eleventh conference on Uncertainty in artificial intelligence*, pages 499–506, 1995.

- [67] Viacheslav Surkov, Chris Wendler, Mikhail Terekhov, Justin Deschenaux, Robert West, and Caglar Gulcehre. Unpacking sdxl turbo: Interpreting text-to-image models with sparse autoencoders. arXiv preprint arXiv:2410.22366, 2024.
- [68] Gemma Team. Gemma. Kaggle, 2024.
- [69] Yu-Lin Tsai, Chia-Yi Hsu, Chulin Xie, Chih-Hsun Lin, Jia-You Chen, Bo Li, Pin-Yu Chen, Chia-Mu Yu, and Chun-Ying Huang. Ring-a-bell! how reliable are concept removal methods for diffusion models? *arXiv preprint arXiv:2310.10012*, 2023.
- [70] 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. Advances in neural information processing systems, 34:16451–16467, 2021.
- [71] Andrey Voynov and Artem Babenko. Unsupervised discovery of interpretable directions in the gan latent space. In *International conference on machine learning*, pages 9786–9796. PMLR, 2020.
- [72] Feng Xie, Biwei Huang, Zhengming Chen, Yangbo He, Zhi Geng, and Kun Zhang. Identification of linear non-gaussian latent hierarchical structure. In *International Conference on Machine Learning*, pages 24370–24387. PMLR, 2022.
- [73] Danru Xu, Dingling Yao, Sébastien Lachapelle, Perouz Taslakian, Julius Von Kügelgen, Francesco Locatello, and Sara Magliacane. A sparsity principle for partially observable causal representation learning. *arXiv preprint arXiv:2403.08335*, 2024.
- [74] Dingling Yao, Danru Xu, Sebastien Lachapelle, Sara Magliacane, Perouz Taslakian, Georg Martius, Julius von Kügelgen, and Francesco Locatello. Multi-view causal representation learning with partial observability. In *The Twelfth International Conference on Learning Rep*resentations, 2023.
- [75] Jaehong Yoon, Shoubin Yu, Vaidehi Patil, Huaxiu Yao, and Mohit Bansal. Safree: Training-free and adaptive guard for safe text-to-image and video generation. *arXiv preprint arXiv:2410.12761*, 2024.
- [76] Jiji Zhang. On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias. *Artificial Intelligence*, 172(16-17):1873–1896, 2008.
- [77] Kun Zhang, Shaoan Xie, Ignavier Ng, and Yujia Zheng. Causal representation learning from multiple distributions: A general setting. In *Forty-first International Conference on Machine Learning*, 2024.
- [78] Nevin L Zhang. Hierarchical latent class models for cluster analysis. *The Journal of Machine Learning Research*, 5:697–723, 2004.
- [79] Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and Sijia Liu. To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now. In *European Conference on Computer Vision*, pages 385–403. Springer, 2024.
- [80] Yujia Zheng, Ignavier Ng, and Kun Zhang. On the identifiability of nonlinear ica: Sparsity and beyond. *arXiv preprint arXiv:2206.07751*, 2022.
- [81] Yujia Zheng, Zeyu Tang, Yiwen Qiu, Bernhard Schölkopf, and Kun Zhang. Detecting and identifying selection structure in sequential data. In *International Conference on Machine Learning*, pages 61498–61525. PMLR, 2024.

Appendix for "Beyond the Black Box: Identifiable Interpretation and Control in Generative Models via Causal Minimality"

A Formulation, Theory, and Experiments for Language Models

A.1 Formulation for Text Generation

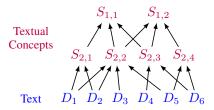


Figure 4: A textual concept graph. We denote text as **D** and discrete textual concepts as **S**. High-level concepts function as selection variables of low-level variables.

Textual concepts are $\mathbf{S} := [\mathbf{S}_1, \cdots, \mathbf{S}_{L_{\mathrm{T}}}]$, where L_{T} is the number of textual hierarchical levels and $\mathbf{S}_l \in \Omega_l \subset \mathbb{N}^{n(\mathbf{S}_l)}$ are concepts at level l.

$$\mathbf{S}_1 \sim \mathbb{P}\left(\mathbf{S}_1\right), \quad \mathbf{S}_l \sim \mathbb{P}\left(\mathbf{S}_l | \mathbf{S}_{l-1}\right), \quad l \in \{2, \dots, L_{\mathrm{T}} + 1\},$$
 where we denote $\mathbf{Z}_0 := \mathbf{D}, \mathbf{Z}_{L_{\mathrm{V}}+1} := \mathbf{X}$, and $\mathbf{S}_{L_{\mathrm{T}}+1} := \mathbf{D}$.

Connections to autoregressive language models. An autoregressive language model can be seen as learning an "encoder" that maps a sequence of input tokens $(D_{1:t})$ to an internal state $\hat{\mathbf{S}}_l$ (e.g., activations within transformer layers). This internal state $\hat{\mathbf{S}}_l$ then informs the "decoder" to predict the subsequent token. For optimal prediction, this learned representation $\hat{\mathbf{S}}_l$ should ideally capture the information of the true concept \mathbf{S}_l that *d-separates* the input tokens $D_{1:t}$ from the next token D_{t+1} . To achieve this d-separation, \mathbf{S}_l should belong to a higher concept level for a larger span of text $D_{1:t}$ (i.e., larger $t \to \text{smaller } l$, see Figure 9). Consequently, broad thematic or narrative structures spanning larger text segments can be compressed into higher-level concepts in our hierarchy (e.g., \mathbf{S}_1), while more localized syntactic or lexical choices correspond to lower-level concepts (e.g., \mathbf{S}_{L_T}).

Intuition on "compression" and higher-level concepts in language models. Our core intuition is that an autoregressive model, at any token position t, compresses the all the preceding sequence (tokens 1 to t) into a representation that is useful for predicting the next token at t+1. In a later position, the model has access to more context and strictly more information. Consequently, the minimality constraint promotes more abstract and compressed representations over the information it has seen. This pressure to compress a growing context naturally gives rise to a hierarchy of concepts. Let's use an example for illustration. When a model reads, "He was secretly buying balloons, sending coded messages to friends, and looking up cake recipes...", it would hold onto this list of disparate actions. The meaning is ambiguous; the model has to keep the details in memory. However, once it has parsed the entire sentence, "He was secretly buying balloons, sending coded messages to friends, and looking up cake recipes – he was getting ready for the surprise party for his sister", the model can now form a high-level concept - a celebratory plan — that organizes all the previous, seemingly random actions into a coherent event. This final concept is more compressed and abstract than the initial list of actions, illustrating the move from detailed memorization to a clear, high-level summary as more context becomes available. In this example, the concepts that exist at later stages of the sequence are not just additions but are fundamentally more abstract, as they synthesize a larger body of information. This aligns directly with our theoretical framework (Condition A.1-iii), where we posit that concepts become more compressed (i.e., have minimal support) as we move up the hierarchy.

A.2 Learning Textual Concepts via State Compression

We now turn to the identification of discrete textual concepts S.

The minimality principle manifests as seeking the most "compressed" representation, namely, achieving minimal support sizes for these discrete concepts while preserving full information.

Condition A.1 (Textual Concept Identification Conditions).

- i Natural Selection: Each selection variable S_l has a support $\operatorname{supp}(S_l)$ that is a proper subset of its potential range if its constituent parts (lower-level variables) were combined randomly. That is, $\operatorname{supp}(S_l) \subsetneq f_{\mathbf{D} \to S_l}(\Omega^{n(\operatorname{Pa}(S_l))})$, where $f_{\mathbf{D} \to S_l}$ is the function from \mathbf{D} to S_l .
- ii **Bottlenecks**: The support size of any concept S_l is strictly smaller than the joint support size of its parents $Pa(S_l)$ in the selection graph.
- iii Minimal Supports: For any S, the condition distribution $\mathbb{P}(\mathbf{D} \setminus \operatorname{Pa}(S)|S=s, \operatorname{HPa}(S)=\tilde{\mathbf{s}})$ is a one-to-one function w.r.t. the argument s.
- iv No-Twins: Distinct latent variables must have distinct sets of adjacent (parent/child) variables.
- v Maximality: The identified latent structure is maximal in the sense that splitting any latent concept variable would violate either the Markov conditions or the No-Twins condition.

Interpreting Condition A.1. Condition A.1-i posits that meaningful text (or textual concepts) occupies a small, structured subset of the vast space of all possible token combinations. We rarely encounter truly random sequences of words in natural language. Conditions A.1-ii and A.1-iii are direct manifestations of causal minimality for discrete concepts. ii implies an information compression moving up the hierarchy—abstract concepts are more succinct. iii demands that each state of a concept s offers unique information about the rest of the text, given its context. Therefore, the representation is most compressed (minimal number of states) and each state contains unique information. Conditions A.1-iv and A.1-v are standard necessary conditions for discrete latent variable model identification [36, 37], precluding redundant or fragmented latent structures.

Theorem A.2 (Textual Concept Identification). Assume the hierarchical process as per (3). Let the true underlying parameters be θ_T . If θ_T satisfies Condition A.1, and an alternative learned model $\hat{\theta}_T$ satisfies Condition A.1-iii, then if both models produce the same observed distribution $\mathbb{P}(\mathbf{D})$, the latent textual concepts \mathbf{S}_l are component-wise identifiable for every level $l \in [L_T]$.

Proof sketch for Theorem A.2. The identification for textual concepts proceeds from the bottom level (tokens, $\mathbf{S}_{L_{\mathrm{T}}}$) upwards to the most abstract concepts (\mathbf{S}_{1}). (1) At each level l+1, we make use of the conditional independence relations that the high-level variable $S_{l,i}$ and its hybrid parents $\mathrm{HPa}(S_{l,i})$ d-separate its pure parents $\mathrm{PPa}(S_{l,i})$ from the other variables $\mathbf{S}_{l}\setminus\{\mathrm{Pa}(S_{l,i})\}$ on level l. This relation allows us to identify subsets of \mathbf{S}_{l+1} that share children on level l [8, 38] and thus reveals the connectivity between variables in \mathbf{S}_{l} and \mathbf{S}_{l+1} . (2) Once the graphical connections are known, we recover the function $\mathrm{Pa}(S_{l,i})\mapsto S_{l,i}$ (i.e., how lower-level concepts combine to form $S_{l,i}$). This is done by merging states of $\mathrm{Pa}(S_{l,i})$ that are predictively equivalent. The "Minimal Supports" (Condition A.1-iii) principle dictates that we choose the function that results in the largest equivalence classes over the parent states (i.e., the most compressed representation for $S_{l,i}$). This ensures that the learned concept $\hat{S}_{l,i}$ has the minimum number of necessary states. (3) This process of structure learning and function recovery is repeated from $\mathbf{S}_{L_{\mathrm{T}}}$ (initially using observed tokens \mathbf{D} as $\mathbf{S}_{L_{\mathrm{T}}+1}$) up to \mathbf{S}_{1} , thereby identifying the entire hierarchy.

Implications for autoregressive language models. Theorem A.2 suggests that by enforcing a minimality regularization for the most compressed representation (Condition A.1-iii), the learned internal states $\hat{\mathbf{S}}$ of a language model can become equivalent to the underlying textual concepts S. SAEs, when applied to transformer activations, can be seen as a practical way to approximate this minimality. By forcing most latent units to be inactive, SAEs force the model to encode information with the minimal active units, which aligns with our theoretical condition for state compression. This result provides a principled justification for the observed interpretability of SAE-derived features and guides our empirical approach in Section A.3 to extract hierarchical textual concept graphs.

A.3 Experiments on Autoregressive Language Models

Implementation. In this section, we present our implementation for analyzing autoregressive language models. We utilize pretrained SAEs [3] for Gemma-2-2b-it [68]. We partition tokens into three parts based on the their positions in their positions in the input sequence. This segmentation reflects the expectation that tokens convey increasingly abstract or high-level information as the

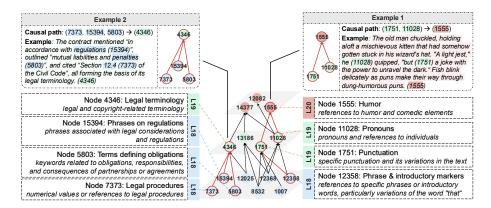


Figure 5: The learned hierarchical concept graph for autoregressive language models. By modeling the hierarchy of concepts based on the token sequence order, we recover a meaningful hierarchical graph. The brown nodes (corresponding to later tokens) capture global, high-level information, while the green nodes (from intermediate tokens) represent more localized, lower-level concepts.

sequence progresses. Finally, we apply causal discovery algorithms to uncover the relationships among features across the different SAEs. More details in Appendix E.

Results. Figure 5 shows a learned hierarchical graph (more in Appendix F). Nodes 1555 and 12082 are mostly activated for final tokens in the sequence, and thus capture high-level semantics. Specifically, node 1555 is associated with the humorous tone, while node 12082 represents the role of the dog. Interestingly, node 11028, derived from intermediate tokens, emerges as a causal factor for both 1555 and 12082. This node encodes pronouns and references to individuals, which play a critical role in shaping both the humor and the characterization of the dog.

B Related Work

Hierarchical models. Complex real-world data distributions frequently exhibit inherent hierarchical structures among their underlying latent variables, a characteristic that has motivated extensive research. Initial explorations primarily focus on continuous latent variables with linear interactions [72, 26, 13, 1]. Other lines of work have centered on discrete latent variables; however, these approaches are often constrained in their applicability to continuous data modalities like images [53, 78, 7, 20, 38]. Furthermore, prevalent latent tree models, which connect variables via a single undirected path [53, 78, 7], risk oversimplifying the multifaceted relationships present in complex systems. More recently, while Park et al. [52] make progress in capturing geometric properties of language model representations using hierarchical models, their work does not address the critical issue of latent variable identification. Kong et al. [39] tackle nonlinear, continuous latent hierarchical models, but their framework, operating under rather opaque functional conditions, falls short of component-wise identifiability, thereby leaving room for concept entanglement. Our work distinctively investigates selection hierarchical models, contending that their structural properties yield a more faithful representation of latent concepts in natural data distributions. In these models, latent variables function as colliders, a significant departure from their role as confounders in the aforementioned prior art. This critical distinction renders existing identification techniques largely inapplicable. To the best of our knowledge, we are the first to provide *component-wise* identifiability for both continuous and discrete hierarchical selection models.

Interpretability for generative models. Despite the remarkable advancements of generative models, their internal mechanisms often remain opaque. This presents a significant challenge to understanding and control. Considerable research has focused on obtaining interpretable features to enable more controllable generation. Early efforts center on analyzing the latent space of generative adversarial networks, e.g., [21, 71, 61]. Recently, sparse autoencoders (SAEs) have gained prominence for interpreting hidden representations, particularly in language models. These studies show that SAEs trained on transformer residual-stream activations can identify latent units corresponding to linguistically meaningful features [11, 28, 18, 51, 62]. These interpretability techniques have also been successfully extended to diffusion models. Recent work [67] reveals interpretable features and specialization across diffusion model blocks. Other work trains SAEs with lightweight classifiers

on diffusion model features [35] or steers generation away from undesirable visual attributes [27]. Distinctly, our work offers a new perspective by framing the concepts learned by generative models within a hierarchical, causal structure. This viewpoint motivates an analytical approach that moves beyond a flat analysis of features. Instead of training a single SAE on aggregated features, our framework suggests dedicating separate SAEs to different stages of the generative process (e.g., distinct diffusion timesteps) to capture concepts at different levels of abstraction. Our theoretical work then provides the formal basis for applying causal discovery across these learned concept layers. This allows for the construction of an explicit hierarchical graph, shifting the focus from an unstructured dictionary of features to an interpretable model of how abstract concepts compose from simpler ones.

C Proofs

C.1 Proof for Theorem 3.3

Lemma C.1 (Base Case Visual Concept Identification). Assume the following data-generating process:

$$\mathbf{C} \sim \mathbb{P}(\mathbf{C}|\mathbf{U}), \mathbf{V} \sim \mathbb{P}(\mathbf{V}), \mathbf{X} := g(\mathbf{C}, \mathbf{V}).$$
 (4)

We have the following conditions.

- *i* Informativeness: The function $g(\cdot)$ is a diffeomorphism.
- ii Smooth Density: The probability density function $p(\mathbf{c}, \mathbf{v}|\mathbf{u})$ is smooth.
- iii Sufficient Variability: At any value \mathbf{c} of \mathbf{C} , there exist $n(\mathbf{C}) + 1$ distinct values of \mathbf{U} , denoted as $\{\mathbf{u}^{(n)}\}_{n=0}^{n(\mathbf{C})}$, such that the vectors $\mathbf{w}(\mathbf{c}, \mathbf{u}^n) \mathbf{w}(\mathbf{c}, \mathbf{u}^0)$ are linearly independent where $\mathbf{w}(\mathbf{c}, \mathbf{u}) = \begin{pmatrix} \frac{\partial \log p(\mathbf{c}|\mathbf{u})}{\partial c_1}, \dots, \frac{\partial \log p(\mathbf{c}|\mathbf{u})}{\partial c_{n(\mathbf{c})}}. \end{pmatrix}$

If a specification θ satisfies i,ii, and iii, another specification $\hat{\theta}$ satisfies i,ii, and they generate matching distribution $\mathbb{P}(\mathbf{X})$, then we can verify that \mathbf{C} and $\hat{\mathbf{C}}$ can be identified up to its subspace.

Proof. Since we have matched distributions, it follows that:

$$p(\mathbf{x}|\mathbf{u}) = \hat{p}(\mathbf{x}|\mathbf{u}). \tag{5}$$

As the generating function g has a smooth inverse (i), we can derive:

$$\begin{split} p(g(\mathbf{c}, \mathbf{v})|\mathbf{u}) &= p(\hat{g}(\hat{\mathbf{c}}, \hat{\mathbf{v}})|\mathbf{u}) \implies \\ p(\mathbf{c}, \mathbf{v}|\mathbf{u}) &|\mathbf{J}_{q^{-1}}| = \hat{p}(g^{-1} \circ \hat{g}(\hat{\mathbf{c}}, \hat{\mathbf{v}})|\mathbf{u}) &|\mathbf{J}_{q^{-1}}| \,. \end{split}$$

Notice that the Jacobian determinant $\left|\mathbf{J}_{g^{-1}}\right| > 0$ because of $g(\cdot)$'s invertibility and let $h := g^{-1} \circ \hat{g}$: $(\hat{\mathbf{c}}, \hat{\mathbf{v}}) \mapsto (\mathbf{c}, \mathbf{v})$ which is smooth and has a smooth inverse thanks to those properties of g and \hat{g} . It follows that

$$p(\mathbf{c}, \mathbf{v}|\mathbf{u}) = \hat{p}(h(\hat{\mathbf{c}}, \hat{\mathbf{v}})|\mathbf{u}) \Longrightarrow p(\mathbf{c}, \mathbf{v}|\mathbf{u}) = \hat{p}(\hat{\mathbf{c}}, \hat{\mathbf{v}}|\mathbf{u}) |\mathbf{J}_{h^{-1}}|.$$

The independence relation in the generating process implies that

$$\log p(\mathbf{c}|\mathbf{u}) + \sum_{i \in [n(\mathbf{v})]} \log p(V_i) = \log \hat{p}(\hat{\mathbf{c}}|\mathbf{u}) + \sum_{i \in [n(\hat{\mathbf{v}})]} \log \hat{p}(\hat{V}_i) + \log |\mathbf{J}_{h^{-1}}|.$$
 (6)

For any realization \mathbf{u}^0 , we subtract (6) at any $\mathbf{u} \neq \mathbf{u}^0$ with that at \mathbf{u}^0 :

$$\log p(\mathbf{c}|\mathbf{u}) - \log p(\mathbf{c}|\mathbf{u}^0) = \log \hat{p}(\hat{\mathbf{c}}|\mathbf{u}) - \log \hat{p}(\hat{\mathbf{c}}|\mathbf{u}^0). \tag{7}$$

Taking derivative w.r.t. \hat{v}_i for $j \in [n(\hat{\mathbf{v}})]$ yields:

$$\sum_{i \in [n(\mathbf{c})]} \frac{\partial}{\partial c_i} (\log p(\mathbf{c}|\mathbf{u}) - \log p(\mathbf{c}|\mathbf{u}^0)) \cdot \frac{\partial c_i}{\partial \hat{v}_j} = 0.$$
 (8)

The left-hand side zeros out because $\hat{\mathbf{c}}$ is not a function of $\hat{\mathbf{v}}$.

Condition iii ensures the existence of at least $n(\mathbf{c})$ such equations with $\mathbf{u}^1, \dots, \mathbf{u}^{n(\mathbf{c})}$ that are linearly independent, constituting a full-rank linear system. Since the choice of $j \in [\mathbf{v}]$ is arbitrary. It follows that

$$\frac{\partial c_i}{\partial \hat{v}_j} = 0, \forall i \in [n(\mathbf{c})], j \in [n(\mathbf{v})]. \tag{9}$$

Therefore, the Jacobian matrix J_h is of the following structure:

$$\mathbf{J}_{h} = \begin{bmatrix} \frac{\partial \mathbf{v}}{\partial \hat{\mathbf{v}}} & \frac{\partial \mathbf{v}}{\partial \hat{\mathbf{c}}} \\ \frac{\partial \mathbf{c}}{\partial \hat{\mathbf{v}}} & \frac{\partial \mathbf{c}}{\partial \hat{\mathbf{c}}} \end{bmatrix}$$
(10)

(9) suggests that the block $\frac{\partial \mathbf{c}}{\partial \hat{\mathbf{v}}} = 0$. Since \mathbf{J}_h is full-rank, we can deduce that $\frac{\partial \mathbf{c}}{\partial \hat{\mathbf{c}}}$ must have full row-rank and $n(\mathbf{c}) \leq n(\hat{\mathbf{c}})$. The sparsity constraint in (2) further implies that $n(\mathbf{c}) = n(\hat{\mathbf{c}})$. That is, we can correctly identify the dimensionality of the changing subspace \mathbf{c} . Moreover, since \mathbf{J}_h is full-rank and the block $\frac{\partial \mathbf{c}}{\partial \hat{\mathbf{v}}}$ is zero, we can derive that the corresponding block $\frac{\partial \hat{\mathbf{c}}}{\partial \mathbf{v}}$ in its inverse matrix $\mathbf{J}_{h^{-1}}$ is also zero. Therefore, there exists an invertible map $\hat{\mathbf{c}} \mapsto \mathbf{c}$, which concludes the proof.

Lemma C.2 (Determining Intersection Cardinality from Union Cardinalities). Let $\mathcal{A} = \{A_1, A_2, \dots, A_n\}$ be a finite collection of finite sets. If for any non-empty subset of indices $K \subseteq \{1, 2, \dots, n\}$, the cardinality of the union $|\bigcup_{k \in K} A_k|$ is known, then for any non-empty subset of indices $S \subseteq \{1, 2, \dots, n\}$, the cardinality of the intersection $|\bigcap_{s \in S} A_s|$ can be determined.

Proof. We proceed by induction on the size of the set of indices S, denoted by |S|, for which we want to determine the intersection cardinality.

Base Case: |S| = 1. Let $S = \{i\}$ for some $i \in \{1, 2, ..., n\}$. We aim to determine the cardinality $\left| \bigcap_{s \in S} A_s \right| = |A_i|$. The union of a single set A_i is simply A_i itself. That is, $A_i = \bigcup_{k \in \{i\}} A_k$. By the premise of the theorem, the cardinality $\left| \bigcup_{k \in \{i\}} A_k \right|$ is known. Therefore, $|A_i|$ is known. The base case holds.

Inductive Hypothesis: Assume that for some integer $m \ge 1$, the cardinality of any intersection of j sets, $\left|\bigcap_{j \in J} A_j\right|$, can be determined from the known union cardinalities for all non-empty index sets J such that $1 \le |J| \le m$.

Inductive Step: We want to show that the cardinality of any intersection of m+1 sets can be determined. Let S_{m+1} be an arbitrary non-empty subset of indices from $\{1,2,\ldots,n\}$ such that $|S_{m+1}|=m+1$. Our goal is to determine $\left|\bigcap_{s\in S_{m+1}}A_s\right|$.

Consider the Principle of Inclusion-Exclusion (PIE) applied to the union of the sets whose indices are in S_{m+1} :

$$\left| \bigcup_{s \in S_{m+1}} A_s \right| = \sum_{\emptyset \neq K \subseteq S_{m+1}} (-1)^{|K|-1} \left| \bigcap_{k \in K} A_k \right|$$

This sum runs over all non-empty subsets K of S_{m+1} . We can separate the term where $K = S_{m+1}$ (which corresponds to the intersection of all m+1 sets) from the other terms in the sum:

$$\left| \bigcup_{s \in S_{m+1}} A_s \right| = \left(\sum_{\emptyset \neq K \subset S_{m+1}} (-1)^{|K|-1} \left| \bigcap_{k \in K} A_k \right| \right) + (-1)^{|S_{m+1}|-1} \left| \bigcap_{s \in S_{m+1}} A_s \right|$$

Here, the sum is now over all non-empty *proper* subsets K of S_{m+1} . We can rearrange this equation to solve for the term $\left|\bigcap_{s\in S_{m+1}}A_s\right|$:

$$(-1)^{|S_{m+1}|-1} \left| \bigcap_{s \in S_{m+1}} A_s \right| = \left| \bigcup_{s \in S_{m+1}} A_s \right| - \sum_{\emptyset \neq K \subset S_{m+1}} (-1)^{|K|-1} \left| \bigcap_{k \in K} A_k \right|$$

Multiplying both sides by $(-1)^{|S_{m+1}|-1}$ (noting that $((-1)^{|S_{m+1}|-1})^2 = 1$):

$$\left| \bigcap_{s \in S_{m+1}} A_s \right| = (-1)^{|S_{m+1}|-1} \left(\left| \bigcup_{s \in S_{m+1}} A_s \right| - \sum_{\emptyset \neq K \subset S_{m+1}} (-1)^{|K|-1} \left| \bigcap_{k \in K} A_k \right| \right)$$

Let us analyze the terms on the right-hand side of this equation:

- 1. The factor $(-1)^{|S_{m+1}|-1}$ is a known sign, since $|S_{m+1}| = m+1$.
- 2. The term $\left|\bigcup_{s\in S_{m+1}}A_s\right|$ is the cardinality of a union of m+1 sets. Since S_{m+1} is a non-empty subset of indices, this value is known by the premise of the theorem.
- 3. Consider the sum $\sum_{\emptyset
 eq K \subset S_{m+1}} (-1)^{|K|-1} \left| \bigcap_{k \in K} A_k \right|$. Each K in this summation is a nonempty proper subset of S_{m+1} . Therefore, the size of each such K satisfies $1 \leq |K| \leq m$. By the Inductive Hypothesis, for any such K (i.e., for any intersection of j sets where $1 \leq j \leq m$), the cardinality $\left| \bigcap_{k \in K} A_k \right|$ can be determined from the known union cardinalities. Consequently, every term in this summation, including its sign factor $(-1)^{|K|-1}$, is determinable.

Since all components on the right-hand side of the equation are known or can be determined based on the theorem's premise and the inductive hypothesis, the value of $\left|\bigcap_{s\in S_{m+1}} A_s\right|$ can be determined.

In conclusion, by the principle of mathematical induction, for any non-empty subset of indices $S \subseteq \{1, 2, \dots, n\}$, the cardinality of the intersection $\left|\bigcap_{s \in S} A_s\right|$ can be determined if the cardinality of any union $\left|\bigcup_{k \in K} A_k\right|$ (for any non-empty $K \subseteq \{1, 2, \dots, n\}$) is known.

Lemma C.3 (Intersection Block Identification [40]). We assume the following data-generating process:

$$[\mathbf{v}_1, \mathbf{v}_2] = g(\mathbf{c}, \mathbf{s}_1, \mathbf{s}_2),\tag{11}$$

$$\mathbf{v}_1 = g_1(\mathbf{c}, \mathbf{s}_1),\tag{12}$$

$$\mathbf{v}_2 = g_2(\mathbf{c}, \mathbf{s}_2),\tag{13}$$

where $\mathbf{c} \in \mathcal{C} \subset \mathbb{R}^{d_c}$, $\mathbf{s}_1 \in \mathcal{S} \subset \mathbb{R}^{d_{s_1}}$, and $\mathbf{s}_2 \in \mathcal{S}_2 \subset \mathbb{R}^{d_{s_2}}$. Both g_1 and g_2 are smooth and have non-singular Jacobian matrices almost everywhere, and g_1 is invertible. If $g_1 : \mathcal{Z} \to \mathcal{V}_1$ and $g_2 : \mathcal{Z} \to \mathcal{V}_2$ assume the generating process of the true model (g_1, g_2) and match the joint distribution $p_{\mathbf{v}_1, \mathbf{v}_2}$, then there is a one-to-one mapping between the estimate $\hat{\mathbf{c}}$ and the ground truth \mathbf{c} over $\mathcal{C} \times \mathcal{S} \times \mathcal{S}$, that is, \mathbf{c} is block-identifiable.

Lemma C.4 (One-level Visual Concept Identification). Assume the process for visual concepts in (1) with $L_{\rm V}=1$. If a model specification $\theta_{\rm V}$ satisfies Condition 3.2, and an alternative specification $\hat{\theta}_{\rm V}$ satisfies Conditions 3.2-i and 3.2-ii, along with a sparsity constraint such that for corresponding \hat{Z} and Z:

$$n(\operatorname{Pa}(\hat{Z})) \le n(\operatorname{Pa}(Z)),$$
 (14)

then, if both models $\theta_{\mathbf{V}}$ and $\hat{\theta}_{\mathbf{V}}$ generate the same observed data distribution $\mathbb{P}(\mathbf{X})$, the latent visual concepts \mathbf{Z}_1 are component-wise identifiable for every level.

Proof. For notational convenience, we denote \mathbb{Z}_1 as \mathbb{S} and \mathbb{D} as \mathbb{U} in this proof. This proof consists of two steps. In step one, we identify the connectivity between U and S variables. In step two, we further show the identifiability of the blocks resulting from intersecting the parent sets $\operatorname{Pa}(U)$ of multiple U variables.

Step 1: connectivity identification. Since we have access to the joint distribution $\mathbb{P}(S, U)$, we can derive conditional distributions $\mathbb{P}(S|\{U_i\}_{i\in\mathcal{H}})$ for any index subset $\mathcal{H}\subseteq[n(U)]$. By Lemma C.1, we can identify the dimensionality of the set of variables S that are connected to *any* variable in $\{U_i\}_{i\in\mathcal{H}}$ for any $\mathcal{H}\subseteq[n(U)]$. Lemma C.2 implies that we can identify the dimensionality of the set of variables S that are connected to *all* variables in $\{U_i\}_{i\in\mathcal{H}}$ for any $\mathcal{H}\subseteq[n(U)]$. This information

gives rise to a partition of S components, in which each part is connected to the same set of U variables. Therefore, we have identified the bipartite graph between S and U up to a permutation.

Step 2: intersection block identification. Denote the indices of S variables that are connected to U_i as $\mathcal{I}(i) \subseteq [n(\mathbf{S})]$. We denote the block of S components connected to all variables in $\{U_i\}_{i\in\mathcal{H}}$ as $\mathbf{S}_{\cap_{i\in\mathcal{H}}\mathcal{I}(i)}$ for any $\mathcal{H}\subseteq [n(\mathbf{U})]$. Thanks to Lemma C.1, we can identify the block $\mathbf{S}_{\mathcal{I}(i)}$ connected to the variable U_i for any $i\in[n(\mathbf{U})]$. Lemma C.3 allows us to identify the intersection of any two blocks $\mathbf{S}_{\mathcal{I}(i)\cap\mathcal{I}(j)}$ for $i\neq j$. Therefore, repeated applications of Lemma C.3 leads to the identification of the intersection block $\mathbf{S}_{\cap_{i\in\mathcal{H}}\mathcal{I}(i)}$ for any $\mathcal{H}\subseteq[n(\mathbf{U})]$. This concludes the proof.

Condition 3.2 (Visual Concept Identification Conditions).

- *i Informativeness:* There exists a diffeomorphism $g_l: (\mathbf{Z}_l, \epsilon_l) \mapsto \mathbf{X}$ for $l \in [0, L]$, where ϵ_l denotes independent exogenous variables.
- ii Smooth Density: The probability density function $p(\mathbf{z}_{l+1}|\mathbf{z}_l)$ is smooth for any $l \in [L_V]$.
- iii Sufficient Variability: For each Z and its parents $\tilde{\mathbf{Z}} := \operatorname{Pa}(Z)$, at any value $\tilde{\mathbf{z}}$ of $\tilde{\mathbf{Z}}$, there exist $n(\tilde{\mathbf{Z}}) + 1$ distinct values of Z, denoted as $\{z^{(n)}\}_{n=0}^{n(\tilde{Z})}$, such that the vectors $\mathbf{w}(\tilde{\mathbf{z}}, z^{(n)}) \mathbf{w}(\tilde{\mathbf{z}}, z^{(0)})$ are linearly independent where $\mathbf{w}(\tilde{\mathbf{z}}, z) = \left(\frac{\partial \log p(\tilde{\mathbf{z}}|z)}{\partial \tilde{z}_1}, \dots, \frac{\partial \log p(\tilde{\mathbf{z}}|z)}{\partial \tilde{z}_{n(\tilde{\mathbf{z}})}}\right)$.
- iv Sparse Connectivity (Minimality): For each parent concept \tilde{Z} , there exists a subset of its children $\mathbf{Z} \subseteq \operatorname{Ch}(\tilde{Z})$ such that their only common parent is \tilde{Z} , i.e., $\bigcap_{Z \in \mathbf{Z}} \operatorname{Pa}(Z) = \{\tilde{Z}\}$.

Theorem 3.3 (Visual Concept Identification). Assume the process for visual concepts in (1). If a model specification θ_V satisfies Condition 3.2, and an alternative specification $\hat{\theta}_V$ satisfies Conditions 3.2-i and 3.2-ii, along with a sparsity constraint such that for corresponding \hat{Z} and Z:

$$n(\operatorname{Pa}(\hat{Z})) \le n(\operatorname{Pa}(Z)),$$
 (2)

then, if both models θ_V and $\hat{\theta}_V$ generate the same observed data distribution $\mathbb{P}(\mathbf{X})$, the latent visual concepts \mathbf{Z}_l are component-wise identifiable for every level $l \in [L_V]$.

Proof. By Lemma C.4, we can identify the set of variables \mathbf{Z}_1 that are directly connected to the text variables \mathbf{D} and their causal graph. Treating the identified \mathbf{Z}_1 as the \mathbf{U} in Lemma C.4, we can further identify \mathbf{Z}_2 . Repeating this procedure yields the identifiability of the entire model.

C.2 Proof for Theorem A.2

Definition C.5 (Non-negative Rank). The non-negative rank of a non-negative matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is equal to the smallest number p such that there exists a non-negative $m \times p$ -matrix \mathbf{B} and a non-negative $p \times n$ -matrix \mathbf{C} such that $\mathbf{A} = \mathbf{BC}$.

Lemma C.6 (Conditional Independence and Nonnegative Rank [8]). Let $\mathbf{P} \in \mathbb{R}^{m \times n}$ be a bi-variate probability matrix. Then its non-negative rank rank₊(\mathbf{P}) is the smallest non-negative integer p such that \mathbf{P} can be expressed as a convex combination of p rank-one bi-variate probability matrices.

Lemma C.7 (One-level Textual Concept Identification). Assume the hierarchical process as per (3) with $L_{\rm T}=1$. Let the true underlying parameters be $\theta_{\rm T}$. If $\theta_{\rm T}$ satisfies Condition A.1, and an alternative learned model $\hat{\theta}_{\rm T}$ satisfies Condition A.1-iii, then if both models produce the same observed distribution $\mathbb{P}(\mathbf{D})$, the latent textual concepts \mathbf{S}_1 are component-wise identifiable.

Proof. For each observed variable D, we search for the *minimal* set of variables $\mathbf{C} \subseteq (\mathbf{D} \setminus D)$ such that the following conditional independence holds:

$$D \perp \mathbf{D} \underbrace{\backslash (\{D\} \cup \mathbf{C})}_{\mathbf{R}} | (\mathbf{C}, \operatorname{Ch}(D)).$$
 (15)

Note that all D, \mathbf{C} , and \mathbf{R} belong to observed variables, and $\mathrm{Ch}(D)$ is latent. Thanks to Condition A.1-ii and Lemma C.6, we can select \mathbf{C} with which the nonnegative rank of the probability table $\mathbf{T}_{D,\mathbf{D}} \setminus (\{D\} \cup \mathbf{C})|_{\mathbf{C}}$ is strictly smaller than the support size of D.

We argue that such C is the group of variables adjacent to the same variable S at the next level as D. In other words, they are the co-parents of D, CoPa(D).

This is because such \mathbf{C} makes 15 hold and thus $\mathrm{CoPa}(D) \subseteq \mathbf{C}$. Otherwise, there would be open paths passing S that induce dependence between D and $\mathrm{CoPa}(D)$, violating the conditional independence relation in (15). Therefore, the minimality constraint would enforce that $\mathbf{C} = \mathrm{CoPa}(D)$. Repeating this procedure to all $D \in \mathbf{D}$, we can construct \mathbf{S} variables at the next level and the adjacency relations between \mathbf{S} and \mathbf{D} .

We proceed to identify the function $\mathbf{D} \mapsto \mathbf{S}$. We refer to D as a pure parent if D is adjacent to only one variable S in the discovered graph. For each S, we denote its pure parents as \mathbf{D}^S and non-pure parents as $\tilde{\mathbf{D}}^S$. We employ the conditional independence relation $\mathbf{D}^S \perp \mathbf{D} \setminus \operatorname{Pa}(\tilde{\mathbf{D}}^S)|(S, \tilde{\mathbf{D}}^S)$ and Condition A.1-iii to identify the value of S, i.e., the function $f_S := (\mathbf{d}^S, \tilde{\mathbf{d}}^S) \mapsto s$.

We first make use of the conditional independence

$$\mathbf{D}^{S} \perp \mathbf{D} \setminus \operatorname{Pa}(S)|(S, \tilde{\mathbf{D}}^{S})$$
(16)

to merge the states of pure parents \mathbf{D}^S conditioned on the non-pure parents $\tilde{\mathbf{D}}^S$. Specifically, we condition on non-pure parents $\tilde{\mathbf{D}}^S = \tilde{\mathbf{d}}^S$ for any $\tilde{\mathbf{d}}^S$ present in the support. We define an equivalence relation \sim over values of $(\mathbf{D}^S, \tilde{\mathbf{D}}^S)$ where $(\mathbf{d}_1^S, \tilde{\mathbf{d}}^S) \sim (\mathbf{d}_2^S, \tilde{\mathbf{d}}^S)$ iff they give rise to an identical conditional distribution $\mathbb{P}\left(\mathbf{D} \setminus \mathrm{Pa}(S) | \mathbf{D}^S = \mathbf{d}_1^S, \tilde{\mathbf{D}}^S = \tilde{\mathbf{d}}^S\right) = \mathbb{P}\left(\mathbf{D} \setminus \mathrm{Pa}(S) | \mathbf{D}^S = \mathbf{d}_2^S, \tilde{\mathbf{D}}^S = \tilde{\mathbf{d}}^S\right)$.

We further resort to a more global conditional independence by considering $(\mathbf{D}^S, \tilde{\mathbf{D}}^S)$ as a meta-variable and all the children $\mathrm{Ch}(\tilde{\mathbf{D}}^S)$ associated with this meta-variable:

$$(\mathbf{D}^{S}, \tilde{\mathbf{D}}^{S}) \perp \mathbf{D} \setminus \operatorname{Pa}(\operatorname{Ch}(\tilde{\mathbf{D}}^{S})) | (\operatorname{Ch}(\tilde{\mathbf{D}}^{S}), \underbrace{\operatorname{Pa}(\operatorname{Ch}(\tilde{\mathbf{D}}^{S})) \setminus \{\mathbf{D}^{S}, \tilde{\mathbf{D}}^{S}\})}_{:=\tilde{\tilde{\mathbf{D}}}^{S}},$$
(17)

where $(\mathbf{D}^S, \tilde{\mathbf{D}}^S)$ has become a pure parent of the latent variable $\operatorname{Ch}(\tilde{\mathbf{D}}^S)$. We further group values $([\mathbf{d}^S], \tilde{\mathbf{d}}^S)$ following the rule that $([\mathbf{d}^S]_1, \tilde{\mathbf{d}}_1^S) \sim ([\mathbf{d}^S]_2, \tilde{\mathbf{d}}_2^S)$ iff $\mathbb{P}\left(\mathbf{D} \setminus \operatorname{Pa}(\operatorname{Ch}(\tilde{\mathbf{D}}^S))|([\mathbf{D}^S], \tilde{\mathbf{D}}^S) = ([\mathbf{d}^S]_1, \tilde{\mathbf{d}}_1^S), \tilde{\tilde{\mathbf{D}}}^S = \tilde{\tilde{\mathbf{d}}}^S\right) = \mathbb{P}\left(\mathbf{D} \setminus \operatorname{Pa}(\operatorname{Ch}(\tilde{\mathbf{D}}^S))|([\mathbf{D}^S], \tilde{\mathbf{D}}^S) = ([\mathbf{d}^S]_2, \tilde{\mathbf{d}}_2^S), \tilde{\tilde{\mathbf{D}}}^S = \tilde{\tilde{\mathbf{d}}}^S\right)$ for each $\tilde{\tilde{\mathbf{d}}}^S$ on the support. That is, conditioning on any $\tilde{\tilde{\mathbf{d}}}^S$ on the support, $([\mathbf{d}^S]_1, \tilde{\mathbf{d}}_1^S)$ and $([\mathbf{d}^S]_2, \tilde{\mathbf{d}}_2^S)$ cannot be distinguished. Thus, we group them into an equivalence class $[(\mathbf{d}^S, \tilde{\mathbf{d}}^S)]$.

Finally, for each equivalent class $[(\mathbf{d}^S, \tilde{\mathbf{d}}^S)]$, we assign a distinct value \hat{s} . This constitutes a function $\hat{f}_S := (\mathbf{d}^S, \tilde{\mathbf{d}}^S) \mapsto \hat{s}$. Due to the deterministic relation from latent variables and their children, \hat{f}_S is well-defined. We denote the random variable $\hat{S} := \hat{f}_S(\mathbf{D}^S, \tilde{\mathbf{D}}^S)$.

In the following, we show that \hat{S} and S are equivalent up to a bijection. We show this by contradiction. Suppose that there existed (s_0,\hat{s}_0) on their respective support, such that their pre-images partially overlapped $(\mathbf{d}_0^S,\tilde{\mathbf{d}}_0^S)\in \hat{f}_S^{-1}(\hat{s}_0)\cap f_S^{-1}(s_0)$ and $\hat{f}_S^{-1}(\hat{s}_0)\neq f_S^{-1}(s_0)$, where $f_S:(\mathbf{d}^S,\tilde{\mathbf{d}}^S)\mapsto s$ represents the true model. Suppose that $f_S^{-1}(s_0)$ missed some elements in $\hat{f}_S^{-1}(\hat{s}_0)$, i.e., $\exists (\mathbf{d}_1^S,\tilde{\mathbf{d}}_1^S)\in \hat{f}_S^{-1}(\hat{s}_0)\setminus f_S^{-1}(s_0)$. In this case, $(\mathbf{d}_0^S,\tilde{\mathbf{d}}_0^S)$ and $(\mathbf{d}_1^S,\tilde{\mathbf{d}}_1^S)$ would lead to distinct values s_0 and s_1 under model f_S^{-1} . By the construction of \hat{f}_S^{-1} , this would indicate $\mathbb{P}(\mathbf{D}\setminus \mathrm{Pa}(S)|S=s_0)=\mathbb{P}(\mathbf{D}\setminus \mathrm{Pa}(S)|S=s_1)$ and $\mathbb{P}\left(\mathbf{D}\setminus \mathrm{Pa}(\mathrm{Ch}(\tilde{\mathbf{D}}^S))|S=s_0,\tilde{\mathbf{D}}^S=\tilde{\mathbf{d}}^S\right)=\mathbb{P}\left(\mathbf{D}\setminus \mathrm{Pa}(\mathrm{Ch}(\tilde{\mathbf{D}}^S))|S=s_1,\tilde{\mathbf{D}}^S=\tilde{\mathbf{d}}^S\right)$ for each $\tilde{\mathbf{d}}^S$ on the support. Since $s_0\neq s_1$, this violates Condition A.1-iii, giving rise to a contradiction.

Suppose that $f_S^{-1}(s_0)$ contains additional elements, i.e., $\exists (\mathbf{d}_2^S, \tilde{\mathbf{d}}_2^S) \in f_S^{-1}(s_0) \setminus \hat{f}_S^{-1}(\hat{s}_0)$. In this case, $(\mathbf{d}_0^S, \tilde{\mathbf{d}}_0^S)$ and $(\mathbf{d}_2^S, \tilde{\mathbf{d}}_2^S)$ would lead to one value s_0 under model f_S^{-1} . By the construction of \hat{f}_S^{-1} , this would indicate either $\mathbb{P}\left(\mathbf{D} \setminus \mathrm{Pa}(S)|\mathbf{D}^S = \mathbf{d}_0^S, \tilde{\mathbf{D}}^S = \tilde{\mathbf{d}}_0^S\right) \neq \mathbb{P}\left(\mathbf{D} \setminus \mathrm{Pa}(S)|\mathbf{D}^S = \mathbf{d}_2^S, \tilde{\mathbf{D}}^S = \tilde{\mathbf{d}}_2^S\right)$

or
$$\mathbb{P}\left(\mathbf{D} \setminus \operatorname{Pa}(\operatorname{Ch}(\tilde{\mathbf{D}}^S)) | ([\mathbf{D}^S], \tilde{\mathbf{D}}^S) = ([\mathbf{d}^S]_0, \tilde{\mathbf{d}}_0^S), \tilde{\tilde{\mathbf{D}}}^S = \tilde{\tilde{\mathbf{d}}}^S\right) \neq \mathbb{P}\left(\mathbf{D} \setminus \operatorname{Pa}(\operatorname{Ch}(\tilde{\mathbf{D}}^S)) | ([\mathbf{D}^S], \tilde{\mathbf{D}}^S) = ([\mathbf{d}^S]_2, \tilde{\mathbf{d}}_2^S), \tilde{\tilde{\mathbf{D}}}^S = \tilde{\tilde{\mathbf{d}}}^S\right) \text{ for some } \tilde{\tilde{\mathbf{d}}}^S \text{ on the support.}$$

By construction of \hat{f}_S , this would violate conditional independence (16) or (17) which the graphical structure implies, which leads to a contradiction.

Therefore, we have shown that for each pair (s,\hat{s}) on their respective support, their pre-images should be identical as long as they intersect: $\hat{f}_S^{-1}(\hat{s}) \cap f_S^{-1}(s) \neq \emptyset \implies \hat{f}_S^{-1}(\hat{s}) = f_S^{-1}(s)$, which is equivalent to that \hat{S} and S are equivalent up to a bijection.

Condition A.1 (Textual Concept Identification Conditions).

- i Natural Selection: Each selection variable S_l has a support $\operatorname{supp}(S_l)$ that is a proper subset of its potential range if its constituent parts (lower-level variables) were combined randomly. That is, $\operatorname{supp}(S_l) \subsetneq f_{\mathbf{D} \to S_l}(\Omega^{n(\operatorname{Pa}(S_l))})$, where $f_{\mathbf{D} \to S_l}$ is the function from \mathbf{D} to S_l .
- ii **Bottlenecks**: The support size of any concept S_l is strictly smaller than the joint support size of its parents $Pa(S_l)$ in the selection graph.
- iii Minimal Supports: For any S, the condition distribution $\mathbb{P}(\mathbf{D} \setminus \operatorname{Pa}(S)|S=s, \operatorname{HPa}(S)=\tilde{\mathbf{s}})$ is a one-to-one function w.r.t. the argument s.
- iv No-Twins: Distinct latent variables must have distinct sets of adjacent (parent/child) variables.
- v Maximality: The identified latent structure is maximal in the sense that splitting any latent concept variable would violate either the Markov conditions or the No-Twins condition.

Theorem A.2 (Textual Concept Identification). Assume the hierarchical process as per (3). Let the true underlying parameters be θ_T . If θ_T satisfies Condition A.1, and an alternative learned model $\hat{\theta}_T$ satisfies Condition A.1-iii, then if both models produce the same observed distribution $\mathbb{P}(\mathbf{D})$, the latent textual concepts \mathbf{S}_l are component-wise identifiable for every level $l \in [L_T]$.

Proof. By Lemma C.7, we can identify the set of variables S_1 adjacent to D and the bipartite causal graph between these two sets of variables. We then employ the identified S_1 to serve as D in the first step to identify S_2 . Repeating this procedure yields the identifiability of the entire model. \Box

D Key Concept Discussions

The roles and purposes of "Selection-based hierarchy and causality minimality". The selection-based hierarchy and causal minimality are constraints on the natural data distribution (images or text), which is a standard modeling practice in causal representation learning [57]. Specifically, the selection-based hierarchy considers concepts as effects of their constituent parts [81], while causal minimality assumes this underlying causal graph is sparse in a specific way (e.g., Condition 3.2-iv).

"Innate" hierarchical concept graphs. "Innate" refers to the causal structure inherent in the natural data-generating process itself. Latent concepts in the real world interact (e.g., 'eyes' and 'nose' are components of a 'face'), forming a pre-existing causal structure which we refer to as the "innate concept graph."

True latent variables and their verifications. "True latent variables" follow the standard notion in causal representation learning [57]: they are the disentangled, interpretable, semantic factors of the real-world data-generating process (e.g., age, object pose). This is in contrast to a deep learning model's learned features, which are often an entangled, uninterpretable mixture optimized for a specific training objective. Aligning learned features with true latent variables (referred to as "identification") is the central goal, as it enables reliable interpretation (e.g., "this feature is age") and precise control (e.g., "increase this feature to make the face older"). This is a fundamental question that our work addresses through both theoretical guarantees and empirical validation. Our work provides the guarantee that if the data-generating process fulfills the property of causal minimality and our learning objective enforces this (e.g., via sparsity), the model's learned features are provably equivalent to the true latent variables. We then validate this empirically via intervention, a standard practice in causal research [57]. Our experiments (Figure 3 and Figure 5) show that manipulating the theoretically identified features provides semantic control over the generated output, providing evidence that these features are the meaningful causal levers of the generative process.

Validity of the conditions. While assumptions on the unobserved data-generating process may not be validated directly, we have reasoned for the plausibility of our conditions by reflecting on natural properties of real-world data. Beyond standard regularity assumptions like smoothness and variability [32, 33, 29, 30, 77], our key minimality conditions—Sparse Connectivity (Condition 3.2-iv) for vision and Minimal Supports (Condition A.1-ii,iii) for text—are motivated by the observation that concepts typically arise from a sparse set of causes [44, 73, 45, 80, 50] and that language is inherently structured and compressible [60, 10]. Perhaps a more convincing validation is the empirical results. Our experiments provide strong indicative support for these assumptions: by actively enforcing sparsity/compression via SAEs, we successfully extract meaningful concept hierarchies in both vision (Figure 3) and text (Figure 5) that are otherwise dense and not easily interpretable. This success provides support for the usefulness of our overall approach and the validity of our assumptions. We acknowledge that these assumptions, like any in this field, may not hold universally. Fortunately, our strong empirical results suggest they seem effective and plausible for the complex, real-world data we study.

Concept variable interpretation. Our theory proves the existence of a clean, one-to-one mapping between a learned feature and a true latent variable. This guarantee is what makes a principled interpretation possible in the first place. The subsequent step—assigning a human-understandable description to this now-identified concept—is intrinsically a task that requires human validation. This is a fundamental aspect of all interpretability research (perhaps modern vision-language models have the potential to automate this process).

Comparison with recent work [12]. On the technique side, Cywinski et al. [12] feature an elegant concept location technique by utilizing the score function, which could significantly benefit our algorithm. For example, we could employ SAeUron [12] to confirm whether our features at various timesteps match the concept location it identifies. Our causal learning algorithm explicitly learns the inter-connectivity among concepts across hierarchical levels. Thus, to modify a part of a high-level concept, we could focus our scope on only the variables connected to this specific highlevel concept, which lowers the search complexity. In our experiment example, to implement two changes, "replacing the rock with tree stump" and "adding texture to tree stump", SAeUron may need to perform two independent searches across all timesteps and node indices. Our method can help reduce the search space to only the low-level nodes connected to "tree stump". In addition, pinpointing specific diffusion timesteps to intervene on potentially aids in managing undesirable artifacts. Moreover, our explicit concept graph could also give an interpretable, intuitive characterization of the model's knowledge. On the message side, Cywinski et al. [12] propose a novel score function to select the timestep and node index for accurate concept unlearning. Our work's focus is to provide concise and informative theoretical conditions to understand concept learning in both vision and language modalities, with potential applications like concept easing or controllable generation. With this work, we hope the theoretical insights will facilitate the development of refined and dedicated methods in the community.

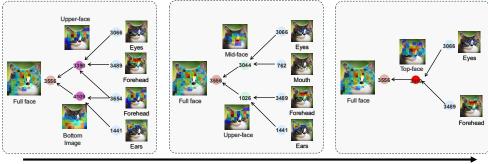
Comparison with recent work [35]. Revelio [35] relies on training a classifier on a specific classification dataset. Revelio trains SAEs and a classifier on a specific dataset (e.g., Caltech-101) to evaluate which features and timesteps are most correlated with class labels. Our work, in contrast, does not involve class labels. Our primary contribution is a hierarchical, causal framework designed to interpret the generative process itself. We apply causal discovery algorithms to discover the causal relationships across different levels of concepts without any class labels. We are able to understand how semantic concepts causally relate to one another across different levels of abstraction to form a coherent output (e.g., how "ear" and "mouth" features causally contribute to a "cat face"). Moreover, the work [35] neither performs interventions nor analyzes the compositional structure of generation, which are the central themes of our paper.

E Implementation Details

We present the diagram of our method in Fig.9.

Computing resources. We use one L40 GPU for training the SAEs and a standard MacBook Pro with an M1 chip for causal discovery. Training one SAE takes around 8 hours.

Vision experiments. For the diffusion sampling process, we utilize the sde-dpmsolver++ [48] sampler, which adds stochasticity between successive steps. We train the K-sparse autoencoder



Sparsity Increasing

Figure 6: Ablation studies on the sparsity constraint. We control feature sparsity at timestep 500. Without enforcing sparsity, the resulting concepts tend to be dense, and the features are less interpretable. Conversely, higher sparsity leads to a more interpretable, sparser graph. However, when sparsity becomes too high, the resulting graph may become overly sparse and fail to adequately capture the generation of the cat face.

using a latent dimension of 5120, a batch size of 4096, and the Adam optimizer with a learning rate of 0.0001, setting K=10. A 10k subset of prompts is selected from LAION-COCO [59]. We then extract SAE features from layers down . 2 . 1, mid . 0, and up . 1 . 0 at timesteps 899,500,100 respectively. We use prompts from the Laion-COCO dataset [59]. The PC algorithm subsequently uses all resulting feature indices for causal discovery. For the sparsity ablation study, we control the top-K value used in the SAE. Specifically, we train additional SAEs with K=4 and K=100 at timestep 500. To evaluate the effect of sparsity (Figure 6), we then perform causal discovery by replacing the SAE features with K=10 with those from the K=4 or K=100 models. Table 1 evaluates the following baselines: SD1.4 [56], ESD [15], SA [22], CA [42], MACE [49], UCE [16], RECE [19], SDID [46], SLD-MAX [58], SLD-STRONG [58], SLD-MEDIUM [58], SD1.4-NegPrompt [56], SAFREE [75], TRASCE [31], and ConceptSteer [34].

LLM experiments. We utilize the pretrained SAEs for gemma-2-2b-it available from Gemma-Scope [68]. To collect features, we use the pile-10k corpus [17]. For each sample, we first exclude padding tokens and divide the remaining meaningful tokens into three sequential segments. The first segment is processed through the SAE at layer 18 to obtain feature indices representing lower-level information. The second segment is passed through the SAE at layer 19 to capture intermediate-level features. The final segment is input to the SAE at layer 20 to extract higher-level features. We then apply the PC algorithm for causal discovery using the feature indices from these three representational levels.

F Additional Empirical Results

Hierarchical causal analysis. Our theoretical framework motivates an empirical analysis that differs from standard interpretability approaches. To empirically instantiate our theory, we first analyze feature representations in Stable Diffusion (SD) 1.4 at distinct timesteps (899, 500, and 100) that our theory posits correspond to different levels of abstraction. Our approach involves two key steps: 1) Level-specific concept learning: In line with our hierarchical model, we train a separate K-sparse autoencoder for each timestep's feature set, allowing us to learn concept dictionaries specific to each level. 2) Cross-level causal discovery: subsequently, to map the relationships predicted by our theory, we apply a causal discovery algorithm (PC [66]) on the learned sparse features across these levels to construct the hierarchical concept graph.

Benefits. This hierarchical perspective provides two main benefits. First, it enables compositional editing. For a complex object like "a textured tree stump", our analysis can distinguish the "stump" (a mid-level concept) from its "texture" (a low-level one), allowing for independent steering. This is a fine-grained control challenging for non-hierarchical methods that tend to learn entangled features (see Table 8). Second, it allows for targeted intervention. By identifying a concept's level, we can inject a steered feature back into the diffusion process only at its corresponding timestep, which helps in reducing the unwanted artifacts that can arise from applying steering globally across all timesteps (see Figure 7). More details in Appendix E and Figure 9.

Metric	SD 1.4	SD1.4 (SAE w/o hier.)	SD1.4 (Ours)
Add tabby pattern – CLIP-I ↓	0.91 ± 0.05	0.83 ± 0.07	0.93 ± 0.04
Add tabby pattern – CLIP-T ↑	0.27 ± 0.00	$\textbf{0.28} \pm \textbf{0.02}$	$\textbf{0.28} \pm \textbf{0.01}$
Add mountains – CLIP-I ↓	0.84 ± 0.06	0.83 ± 0.04	$\textbf{0.91} \pm \textbf{0.03}$
Add mountains – CLIP-T ↑	$\textbf{0.33} \pm \textbf{0.01}$	0.32 ± 0.01	$\textbf{0.33} \pm \textbf{0.01}$
Replace rock w/ stump – CLIP-I ↓	0.93 ± 0.02	0.95 ± 0.02	$\textbf{0.96} \pm \textbf{0.02}$
Replace rock w/ stump – CLIP-T↑	$\textbf{0.31} \pm \textbf{0.01}$	0.29 ± 0.01	$\textbf{0.31} \pm \textbf{0.01}$

Table 2: Controllable image generation results. Our method achieves the best CLIP-I metric, demonstrating greater fidelity to the input images, while reliably executing the target edits.



Figure 7: **Generated samples with P4D prompts [6].** The Stable Diffusion model is vulnerable to the prompts in the p4d dataset, producing unsafe images. When the hierarchical relationship across timesteps is not considered, negative steering with SAE results in drastic changes to the output. In contrast, our method learns to apply modifications to the nudity feature at a suitable timestep without introducing additional distortions.

Ablation. As established in the theoretical framework, sparsity is crucial for identifiability. To empirically validate this, we visualize the resulting causal graphs under varying levels of sparsity, as shown in Fig. 6 (more in Appendix F). When sparsity is not enforced, the resulting graph becomes overly dense, making it difficult to interpret and diminishing its semantic clarity. Conversely, imposing excessive sparsity leads to an overly pruned graph that lacks sufficient structure to meaningfully explain the generation process, such as in the case of the cat image. These observations highlight the importance of balancing sparsity to preserve interpretability while maintaining explanatory power.

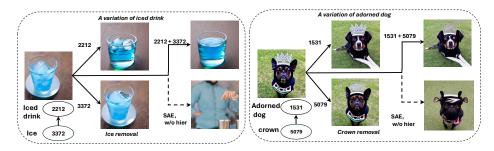


Figure 8: **Examples of multi-level editing** (best viewed with zoom). High-level node 2212 contains all information about the cup, while mid-level node 3372 focuses primarily on the ice cubes. Similarly, high-level node 1531 encompasses all information about the dog (including the crown), and mid-level node 5079 is dedicated to the crown. By modeling hierarchical relationships, we can perform edits that are often difficult to achieve with a single-layer edit. For instance, if we want to generate a variation of the cup while removing the ice cubes, we can apply feature steering on high-level node 2212 to create a new version of the cup, and simultaneously apply negative feature steering on mid-level node 3372 to remove the ice cubes.

F.1 Downstream Tasks

Thanks to our theoretical framework, we can naturally perform a range of image generation and editing tasks, including model unlearning, controllable image generation, and multi-level editing.

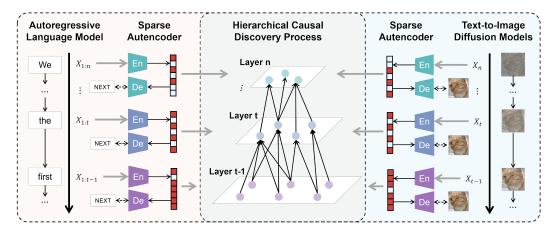


Figure 9: Diagram of our interpretability method. We train SAEs to capture features at different levels (timesteps for diffusion models and token positions for LLMs), and apply causal discovery to construct a hierarchical concept graph.

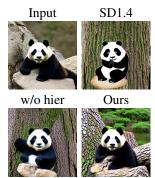


Figure 10: Examples of controllable image generation.

Model unlearning. We provide quantitative results of model unlearning on four benchmark datasets: IP2P [58], three splits of RING-A-BELL [69], P4D [6], and UnlearnDiffATK [79]. These benchmarks focus on removing nudity-related concepts, and we report the accuracy of a pretrained nudity detector. Our method achieves the best results across all benchmarks. In addition, to assess whether our method preserves general text-to-image capability, we apply feature steering on normal prompts from MSCOCO [47]. The 10K results, reflected in low FID and high CLIP scores, demonstrate that our method successfully identifies and removes nudity concepts without affecting unrelated concepts. We also provide results on style removal in the appendix (Table 3) and we achieve superior performance across different metrics and tasks.

Controllable image generation. We also evaluate controllable image generation on three editing tasks: adding tabby patterns to cat faces, adding mountains to landscape images, and replacing rocks with textured tree stumps. As shown in Table 2 and Fig.10, our method achieves superior results compared to both the standard text-guided model and SAE without hierarchical modeling.

Multi-level image editing. A key advantage of the hierarchical concept graph is that it can combine nodes across different levels for fine-grained image editing. In Fig. 8, to obtain a new drink without ice (while preserving the background), we can apply multi-level editing by steering features at both high-level node 2212 and mid-level node 3372 simultaneously. Without such hierarchical relationship modeling, conventional methods struggle to produce this combination, which can result in undesired changes such as the drink being replaced by a person or the dog's background.

More examples for Figure 3. Figure 11 and Figure 12 contain more examples of Figure 3. For example, node 3641 in the SAE at timestep 899 contains comprehensive information about the panda, as illustrated by the heatmap. When feature steering is applied, it results in the generation of a new panda. Meanwhile, nodes 1026 and 511 in the SAE at timestep 500 represent different components of the panda. At a finer level of detail, nodes 3489, 3880, and 451 in the SAE at timestep

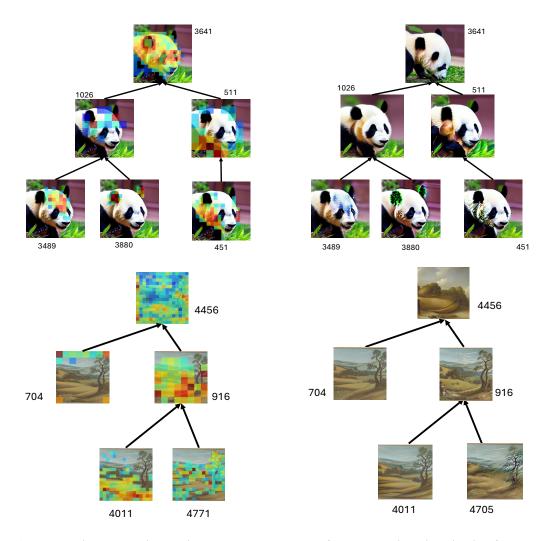


Figure 11: Discovered hierarchical concept graphs and feature steering visualization for text-to-image generation. We can observe that features on the hierarchical model represent a part-whole relation, and steering a feature yields corresponding visual variation (e.g., the panda's ears).

100 capture specific image features. These hierarchical concept graphs effectively illustrate how the panda is generated.

More results for model unlearning In addition to the four benchmark datasets in the main paper, we report results on another commonly used benchmark dataset with two tasks: *Remove Van Gogh* and *Remove Kelly McKernan* in Table.3. We evaluate performance using four metrics: LPIPSe (similarity for prompts with the target style), LPIPSu (similarity for prompts without the style), Acce (how well the target style was removed), and Accu (how well other styles were preserved), with accuracy ratings assessed using GPT-4o. Our method achieves competitive performance across all metrics and tasks.

Understanding the sparsity constraint. Figure 13 and Table 4 contain the ablation study for the sparsity constraint. We can observe that a proper sparsity strength can indeed give rise to desirable interpretability results, while too small and too large sparsity constraints may be harmful in practice. As shown in Table 4, a low sparsity penalty results in visualized maps with significant overlap. On the other hand, applying a strong sparsity penalty leads to low node coverage, indicating that the nodes alone are insufficient to fully explain the generation of the entire image.

More examples for Figure 5. Figure 14 contains more examples for Figure 5. As discussed in the main paper, we divide the tokens into three segments based on their sequence order, with later

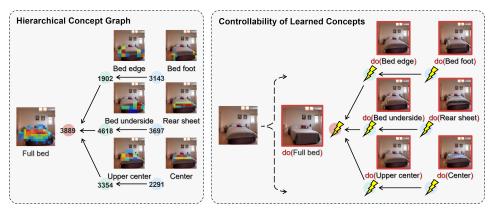


Figure 12: More examples of the learned hierarchical concept graphs for text-to-image models. Under appropriate sparsity and noise conditions, our method successfully recovers meaningful hierarchical structures, where each node encodes distinct semantic concepts. On the right, we demonstrate feature steering, where manipulating individual nodes leads to changes in the output that align with their position in the hierarchy – higher-level nodes produce broader semantic shifts, while lower-level nodes control more fine-grained aspects.

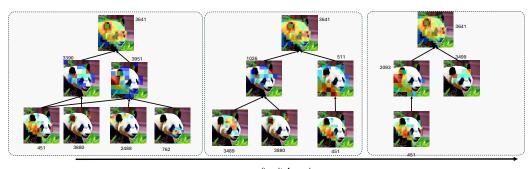
Method	LPIPSe ↑	LPIPSu↓	Acce ↓	Accu ↑	
Task: Remove "Van Gogh"					
SD-v1.4	_	_	0.95	0.95	
CA [42]	0.30	0.13	0.65	0.90	
RECE [19]	0.31	0.08	0.80	0.93	
UCE [16]	0.25	0.05	0.95	0.98	
SLD-Medium [58]	0.21	0.10	0.95	0.91	
SAFREE [75]	0.42	0.31	0.35	0.85	
Ours	0.53	0.26	0.30	0.88	
Task: Remove "Kelly McKernan"					
SD-v1.4	_	_	0.80	0.83	
CA [42]	0.22	0.17	0.50	0.76	
RECE [19]	0.29	0.04	0.55	0.76	
UCE [16]	0.25	0.03	0.80	0.81	
SLD-Medium [58]	0.22	0.18	0.50	0.79	
SAFREE [75]	0.40	0.39	0.40	0.78	
Ours	0.48	0.20	0.35	0.81	

Table 3: **Results on style removal.** We apply negative feature steering to the node to suppress the styles in the image.

tokens expected to encode higher-level information—consistent with the behavior of autoregressive language models. At the highest level, node 11859 represents the "yell mode," characterized by capitalized words conveying a strong tone. The green node 1033, located at an intermediate sequence position, emphasizes importance or intensity—typically a component of the yell mode. At the lowest level, nodes 304, 2009, and 2818 capture various aspects and meanings related to the concept of importance.

-	Overlap ↓	Coverage ↑
K=4	0.108 ± 0.128	26.37 ± 17.24
K=10	0.089 ± 0.079	47.90 ± 12.50
K=100	0.235 ± 0.132	37.46 ± 17.31

Table 4: **Quantitative ablation results.** We generate 100 panda images using different random seeds and visualize the feature heatmaps at timestep 500. We adjust the top-K value in the SAE at timestep 500 to control the level of sparsity. To evaluate, we compute the intersection-over-union (IoU) of intermediate heatmaps to measure concept disentanglement, and the union of all features to assess coverage. IoU reflects how distinctly the intermediate concepts are represented, while coverage in percentage indicates the extent to which the intermediate nodes collectively account for the image generation.



Sparsity Increasing

Figure 13: **Understanding the sparsity constraint.** We adjust the top-K value in the SAE at timestep 500 to control the level of sparsity, effectively modifying the sparsity strength of the SAE at this middle layer. As sparsity decreases, the resulting graph becomes denser, introducing many redundant and semantically irrelevant edges. This reduces the overall interpretability of the concept graph. Conversely, increasing sparsity yields a cleaner, more concise graph. However, if sparsity is too high, it may hinder the formation of a complete and interpretable concept graph necessary for image generation.

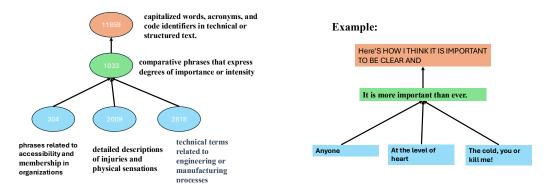


Figure 14: An example of a discovered hierarchical concept graph for autoregressive language modeling. Node 11859 represents a "yell mode," characterized by capitalized words that convey a strong tone. The green node 1033 captures the concept of emphasizing importance or intensity. Blue nodes correspond to lower-level information—for instance, node 304 represents entities mentioned throughout the text.