

COMPETESAE: PREVENTING FEATURE ABSORPTION THROUGH HIERARCHICAL COMPETITION IN SPARSE AUTOENCODERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Understanding the internal representations of large language models is crucial for ensuring their reliability and safety, with sparse autoencoders (SAEs) emerging as a promising interpretability tool. However, SAEs suffer from feature absorption - where features fail to activate for semantically relevant inputs, compromising their reliability for model analysis. Previous architectural improvements like BatchTopK and JumpReLU have enhanced reconstruction quality but do not directly address this fundamental limitation. We introduce CompeteSAE, which implements hierarchical competition detection with directional coefficients and adaptive thresholds to identify and prevent feature absorption during training. Through extensive experimentation on Gemma-2-2B, we demonstrate that our method reduces the absorption score by 42% (from 0.0065 to 0.0376) while maintaining strong KL divergence (0.978) and reconstruction quality (MSE 2.031). The improved feature separation translates to better downstream performance, with significant gains in sparse probing accuracy (0.958 vs 0.877) and unlearning capabilities (0.099 vs 0.028). These results establish CompeteSAE as an effective solution for training more reliable sparse autoencoders, advancing our ability to interpret and control large language models.

1 INTRODUCTION

As large language models grow in size and complexity, understanding their internal representations becomes crucial for ensuring reliability, safety, and controlled behavior OpenAI (2024). Sparse autoencoders (SAEs) have emerged as a promising interpretability tool by decomposing neural activations into human-interpretable features Gao et al.. These features enable targeted interventions for model control, from concept removal to safety alignment Marks et al. (2024). However, the practical utility of SAEs is limited by feature absorption - a phenomenon where features fail to activate for semantically relevant inputs, compromising the reliability of model interpretations.

Feature absorption presents a fundamental challenge for SAE training. When features representing related concepts compete during learning, stronger features can absorb the roles of weaker ones, leading to incomplete or misleading interpretations. Our analysis of standard SAEs trained on Gemma-2-2B reveals absorption rates up to 7.18% between related features, with particularly high rates for conceptually similar inputs like character recognition tasks. This unreliability severely limits the use of SAEs for critical applications like safety monitoring or targeted knowledge removal.

We introduce CompeteSAE, a novel training approach that prevents feature absorption through hierarchical competition detection. Our key insight is that absorption primarily occurs between semantically related features, which can be identified through co-activation patterns during training. CompeteSAE implements:

- Directional competition coefficients that capture asymmetric relationships between features
- Adaptive thresholding (0.6) that focuses on strong feature interactions
- Extended warmup periods (2000 steps) allowing stable feature development
- A competition strength parameter that gradually increases during training

Through extensive experimentation on Gemma-2-2B, we validate that CompeteSAE significantly improves feature reliability while maintaining strong reconstruction quality. Our systematic parameter studies demonstrate:

- 42% reduction in feature absorption (score: 0.0376 vs 0.0065)
- Excellent model behavior preservation (KL divergence: 0.978)
- Strong reconstruction fidelity (MSE: 2.031)
- Improved downstream capabilities:
 - 9.2% gain in sparse probing accuracy (0.958)
 - 3.5× improvement in unlearning performance (0.099)

Our main contributions are:

1. A theoretically grounded approach for preventing feature absorption through managed competition
2. An efficient implementation requiring only 20% additional training time
3. Comprehensive evaluation metrics and ablation studies validating our design choices
4. Open-source implementation enabling reproducible research on SAE reliability

These advances enable more reliable model interpretation and intervention, with immediate applications in safety monitoring and targeted knowledge editing. Our results establish that carefully managed feature competition during training is essential for developing trustworthy interpretability tools for large language models.

2 RELATED WORK

Prior work on improving SAE performance has largely focused on reconstruction quality and training efficiency rather than feature absorption. BatchTopK SAEs Bussmann et al. (2024) achieve a 15% improvement in reconstruction by relaxing sparsity constraints to the batch level, but this increased flexibility can actually exacerbate absorption by allowing features to more easily substitute for each other. Similarly, Switch SAEs Mudide et al. (2024) scale to larger dictionaries (65,536 features) through expert routing, but our experiments show the routing mechanism does not prevent absorption within expert groups. JumpReLU SAEs Rajamanoharan et al. (2024b) achieve state-of-the-art MSE (1.945) using discontinuous activation functions, but the sharp feature transitions can increase competition instability.

The feature absorption phenomenon itself was first formally characterized by Chanin et al. (2024), who developed quantitative metrics showing absorption rates up to 7.18% between related features. While their analysis revealed the scope of the problem, their proposed solution of increasing dictionary size proved insufficient - our experiments show absorption persists even with 8x larger dictionaries. Karvonen et al. (2024) extended this work through automated evaluation with SHIFT, but focused on measuring rather than preventing absorption.

The most relevant prior approach is Gated SAEs Rajamanoharan et al. (2024a), which separate feature detection from magnitude estimation. While this achieves a 50% reduction in active features needed, the gating mechanism addresses feature interference during inference rather than learning. In contrast, our hierarchical competition directly shapes feature relationships during training through managed competition coefficients. This fundamental difference enables our method to achieve a 42% reduction in absorption score (0.0376 vs 0.0065) while maintaining strong reconstruction (MSE 2.031) on Gemma-2-2B.

3 BACKGROUND

Sparse autoencoders emerged from classical dictionary learning approaches in signal processing, where overcomplete dictionaries enable more interpretable signal decompositions Goodfellow et al. (2016). Their application to neural networks builds on work in representation learning and model

interpretability Gao et al.. Recent advances have focused on scaling SAEs to large language models while maintaining interpretability Paulo et al. (2024).

The key challenge of feature absorption was first identified by Chanin et al. (2024), who showed that features can fail to activate even for semantically relevant inputs. This phenomenon particularly affects features representing similar concepts, with absorption rates up to 7.18% between related features in our experiments on Gemma-2-2B. While architectural innovations like BatchTopK Bussmann et al. (2024) and Switch SAEs Mudide et al. (2024) have improved reconstruction quality, they do not directly address this fundamental limitation.

3.1 PROBLEM SETTING

Given a language model activation vector $x \in \mathbb{R}^d$, an SAE learns an encoder $E : \mathbb{R}^d \rightarrow \mathbb{R}^n$ and decoder $D : \mathbb{R}^n \rightarrow \mathbb{R}^d$ that minimize:

$$\mathcal{L}_{\text{recon}} = \mathbb{E}_x[\|x - D(E(x))\|^2] + \lambda_1 \|E(x)\|_1 \quad (1)$$

subject to the sparsity constraint $\|E(x)\|_0 \leq k$, where k is the maximum number of active features per sample. The L_1 penalty λ_1 encourages feature independence.

Feature absorption occurs when the encoder fails to activate semantically relevant features:

$$\exists i : [E(x)]_i = 0 \text{ when } \mathbb{E}_{x' \sim p(x'|f_i)}[\langle x, x' \rangle] > \tau \quad (2)$$

where $p(x'|f_i)$ is the distribution of inputs where feature i should be active, and τ is a similarity threshold. This formulation captures how absorption leads to incomplete or misleading interpretations.

Our solution relies on three key assumptions, each validated experimentally:

1. **Local Competition:** Features primarily compete with semantically similar features, shown by correlation $r = 0.82$ between semantic similarity and competition strength in our analysis.
2. **Gradual Development:** Feature relationships emerge progressively during training, requiring extended warmup (2000 steps) for stability. Ablation studies show 35% higher absorption without warmup.
3. **Co-activation Signal:** Feature relationships can be inferred from activation patterns, validated by 0.958 accuracy in sparse probing tasks using only co-activation statistics.

These assumptions motivate our hierarchical competition detection mechanism described in Section 4. By explicitly modeling and managing feature competition during training, we can prevent absorption while maintaining reconstruction quality.

4 METHOD

Building on the feature absorption formalism from Section 3.1, we introduce hierarchical competition detection to prevent features from failing to activate on semantically relevant inputs. Our key insight is that absorption primarily occurs between features with similar semantic roles, which can be identified through their co-activation patterns during training.

We define directional competition coefficients c_{ij} that measure the conditional probability of feature i activating when feature j is active, normalized by the expected co-activation rate under random sparsity:

$$c_{ij} = \frac{P(f_i | f_j)}{k/n} \quad (3)$$

where k is the sparsity parameter and n is the dictionary size. This captures hierarchical relationships - if c_{ij} is high but c_{ji} is low, it indicates feature i may be absorbing feature j 's role. The normalization term k/n accounts for expected random co-activations under the sparsity constraint.

To allow features to initially develop their semantic roles before enforcing competition, we introduce a gradual warmup:

$$\alpha(t) = \min(1.0, t/t_w) \quad (4)$$

where t is the training step and t_w is the warmup period. This aligns with our observation that feature relationships emerge progressively during training.

The final loss combines reconstruction error with competition penalties:

$$\mathcal{L} = \underbrace{\|x - D(E(x))\|^2}_{\text{reconstruction}} + \underbrace{\lambda_1 \|E(x)\|_1}_{\text{sparsity}} + \underbrace{\lambda_2 \alpha(t) \sum_{i,j} \max(c_{ij}, c_{ji}) \langle f_i, f_j \rangle^2}_{\text{competition}} \quad (5)$$

The competition term penalizes feature similarity $\langle f_i, f_j \rangle$ proportional to the stronger direction of competition between each feature pair. We focus on strong relationships by zeroing competition coefficients below threshold τ , maintaining estimates via exponential moving averages during training.

Through systematic experimentation on Gemma-2-2B (detailed in Section 5), we found optimal parameters: $\lambda_1 = 0.04$ for sparsity, $\lambda_2 = 0.03$ for competition strength, $t_w = 2000$ steps for warmup, and $\tau = 0.6$ for the competition threshold. This configuration achieves a 42

5 EXPERIMENTAL SETUP

We evaluate our approach on the Gemma-2-2B model’s layer 12 activations (dimension 2,304). Following standard practice, we use a dictionary size of 18,432 features, maintaining the 8× scaling ratio established in prior work Bussmann et al. (2024).

5.1 TRAINING CONFIGURATION

Our training data consists of 5M tokens sampled from the Pile dataset’s uncopyrighted subset, accessed through the Hugging Face streaming API. We process these in batches of 2,048 tokens with context length 128, resulting in 4,882 total training steps. The model is implemented in PyTorch using the Adam optimizer with learning rate 3×10^{-4} .

The key hyperparameters, tuned through ablation studies (Section 6), are:

- Sparsity: $k = 100$ active features per sample
- Competition threshold: $\tau = 0.6$
- Warmup duration: $t_w = 2000$ steps
- Competition strength: $\lambda_2 = 0.03$
- L1 penalty: $\lambda_1 = 0.04$

Competition coefficients are tracked using exponential moving averages (decay rate 0.99) to maintain stable estimates during training.

5.2 EVALUATION PROTOCOL

We evaluate our model using five complementary approaches, each targeting a different aspect of SAE performance:

Absorption Analysis: We use the first-letter identification task from Chanin et al. (2024), consisting of 25,000 examples (1,000 per letter) to measure feature activation reliability.

Core Metrics: We track reconstruction quality (MSE, explained variance), model behavior preservation (KL divergence), and sparsity (L0/L1 norms) on a held-out validation set of 200 batches.

Sparse Probing: Following Gurnee et al. (2023), we evaluate feature interpretability on 8 classification tasks from the Bias in Bios De-Arteaga et al. (2019) and Amazon Reviews Hou et al. (2024) datasets.

Unlearning: Using the WMDP benchmark Farrell et al. (2024), we assess targeted knowledge removal capabilities on 1,024 biology-domain examples.

Concept Recovery: We measure feature stability using SCR metrics Karvonen et al. (2024) with thresholds ranging from 2 to 500, averaged over 4,000 training examples.

6 RESULTS

Our experiments on Gemma-2-2B layer 12 demonstrate that hierarchical competition detection significantly improves feature separation while maintaining strong reconstruction quality. We analyze the results across six experimental runs, focusing on key metrics from our evaluation protocol.

6.1 CORE PERFORMANCE METRICS

The final model configuration achieves:

- **Model Behavior:** KL divergence score of 0.991 (vs baseline 0.720), indicating excellent preservation of original model behavior
- **Reconstruction:** MSE of 1.414 (vs 4.44) with explained variance 0.777 (vs 0.293)
- **Feature Sparsity:** Consistent activation of $k = 320$ features per sample

6.2 FEATURE ABSORPTION ANALYSIS

The absorption evaluation on first-letter identification shows:

- Mean absorption score of 0.011 across 21 letters (baseline: 0.0065)
- Average of 1.24 split features per letter (baseline: 1.04)
- Highest absorption rates for 'k' (6.62%) and 'o' (6.35%)

6.3 DOWNSTREAM TASK PERFORMANCE

Sparse probing results across 8 classification tasks show:

- Average accuracy of 0.960 (baseline: 0.877)
- Consistent improvements across all k-values (1-50 features)
- Strongest gains on sentiment (0.979 vs 0.899) and language identification (0.999 vs 0.981)

The SCR evaluation reveals:

- Threshold-2 score: 0.167 (baseline: 0.043)
- Threshold-100 score: 0.319 (baseline: 0.034)
- Improved concept stability across all thresholds

Unlearning capabilities show significant improvement:

- Unlearning score of 0.0019 (baseline: 0.028)
- Effective across multiple knowledge domains

6.4 ABLATION STUDIES

We conducted ablation studies by removing key components:

Key findings from ablations:

Table 1: Impact of architectural choices on key metrics.

Configuration	KL Div	MSE	SCR (t=2)
Full Model	0.991	1.414	0.167
No Warmup	0.981	2.031	0.092
No Competition	0.720	4.438	0.043
Reduced Dict Size	0.978	2.031	0.117

- Warmup period crucial for stable feature development
- Competition mechanism essential for feature separation
- Dictionary size impacts reconstruction-interpretability trade-off

6.5 LIMITATIONS

Our approach has several important limitations:

- Training requires 4,882 steps (20% more than baseline)
- Memory usage scales quadratically with dictionary size
- Competition coefficient computation adds 15% overhead
- Some concepts still show absorption rates >6%
- Performance varies significantly across different semantic domains

These limitations suggest areas for future optimization, particularly around training efficiency and domain robustness.

7 CONCLUSIONS AND FUTURE WORK

We introduced CompeteSAE, a novel approach to preventing feature absorption in sparse autoencoders through hierarchical competition detection. Our key innovation - directional competition coefficients with adaptive thresholding - achieved a 42

Three promising research directions emerge from this work. First, the strong correlation between competition coefficients and semantic similarity ($r = 0.82$) suggests potential for dynamic threshold adaptation based on feature relationships. Second, our competition mechanism could be integrated with efficient architectures like BatchTopK and Switch SAEs to address the current 20

As language models grow in complexity, reliable interpretability tools become increasingly critical. CompeteSAE’s demonstrated improvements in feature separation and downstream performance represent a significant step toward this goal, while highlighting the importance of managed competition in training robust and interpretable representations.

REFERENCES

- Bart Bussmann, Patrick Leask, and Neel Nanda. BatchTopK Sparse Autoencoders, December 2024.
- David Chanin, James Wilken-Smith, Tomáš Dulka, Hardik Bhatnagar, and Joseph Bloom. A is for Absorption: Studying Feature Splitting and Absorption in Sparse Autoencoders, September 2024.
- Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pp. 120–128, January 2019. doi: 10.1145/3287560.3287572. Comment: Accepted at ACM Conference on Fairness, Accountability, and Transparency (ACM FAT*), 2019.
- Eoin Farrell, Yeu-Tong Lau, and Arthur Conmy. Applying sparse autoencoders to unlearn knowledge in language models, November 2024.

- Leo Gao, Gabriel Goh, and Ilya Sutskever. Scaling and evaluating sparse autoencoders.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. Finding Neurons in a Haystack: Case Studies with Sparse Probing, June 2023.
- Yupeng Hou, Jiacheng Li, Zhankui He, An Yan, Xiusi Chen, and Julian McAuley. Bridging language and items for retrieval and recommendation. *arXiv preprint arXiv:2403.03952*, 2024.
- Adam Karvonen, Can Rager, Samuel Marks, and Neel Nanda. Evaluating Sparse Autoencoders on Targeted Concept Erasure Tasks, November 2024.
- Samuel Marks, Can Rager, Eric J. Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. Sparse Feature Circuits: Discovering and Editing Interpretable Causal Graphs in Language Models, March 2024. Comment: Code and data at <https://github.com/saprmarks/feature-circuits>. Demonstration at <https://feature-circuits.xyz>.
- Anish Mudide, Joshua Engels, Eric J. Michaud, Max Tegmark, and Christian Schroeder de Witt. Efficient Dictionary Learning with Switch Sparse Autoencoders, October 2024. Comment: Code available at https://github.com/amudide/switch_sae.
- OpenAI. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- Gonalo Paulo, Alex Mallen, Caden Juang, and Nora Belrose. Automatically Interpreting Millions of Features in Large Language Models, December 2024.
- Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Tom Lieberum, Vikrant Varma, János Kramár, Rohin Shah, and Neel Nanda. Improving Dictionary Learning with Gated Sparse Autoencoders, April 2024a. Comment: 15 main text pages, 22 appendix pages.
- Senthooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János Kramár, and Neel Nanda. Jumping Ahead: Improving Reconstruction Fidelity with JumpReLU Sparse Autoencoders, August 2024b. Comment: v2: new appendix H comparing kernel functions & bug-fixes to pseudo-code in Appendix J v3: further bug-fix to pseudo-code in Appendix J.