

DYNAMIC FEATURE DISENTANGLEMENT: ADAPTIVE THRESHOLDING FOR INTERPRETABLE 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, existing SAE approaches face a fundamental challenge: they often produce entangled features with inconsistent activation patterns, making it difficult to identify and analyze specific model behaviors. We address this challenge through a novel threshold-based orthogonality approach that dynamically balances feature importance using gradient-based metrics and activation frequencies, while introducing correlation-aware penalty scaling to promote feature separation. Our method achieves this through three key innovations: adaptive masking based on feature statistics, importance-weighted learning rates, and dynamic threshold adjustment. Experiments on the Gemma-2B model demonstrate significant improvements, including a 9.8% reduction in L0 sparsity (76.85 vs baseline 85.21) while maintaining 94% of baseline reconstruction quality (explained variance 0.291 vs 0.309). The approach shows particular strength in semantic concept recognition with peak SCR scores of 0.143 in the optimal threshold range (20-50), and achieves specialized feature detection as evidenced by consistent absorption patterns across letter categories (mean score 0.0101). These results demonstrate that our method successfully balances the competing objectives of sparse feature extraction and faithful reconstruction, providing a more reliable foundation for transformer model interpretation.

1 INTRODUCTION

Understanding the internal mechanisms of large language models (LLMs) is crucial for ensuring their reliability and safety OpenAI (2024). While sparse autoencoders (SAEs) have emerged as a promising tool for model interpretation, they face a fundamental challenge: feature entanglement. Current approaches struggle to maintain consistent activation patterns while preserving model behavior, making it difficult to identify and analyze specific computational patterns. This challenge is particularly acute in transformer architectures, where high-dimensional representations and complex attention mechanisms create intricate feature interactions that resist traditional interpretation methods.

The core technical challenge lies in optimizing three competing objectives: (1) maintaining faithful reconstruction of model activations, (2) achieving sparse and interpretable feature representations, and (3) ensuring consistent feature activation patterns across different inputs. Traditional autoencoder architectures Goodfellow et al. (2016) typically optimize for reconstruction alone, while existing sparse coding approaches Olshausen & Field (1996) lack mechanisms for handling the dynamic, context-dependent nature of transformer features Vaswani et al. (2017). Our baseline experiments on Gemma-2B demonstrate this challenge, achieving high reconstruction quality (explained variance 0.309) but with poor feature separation (L0 sparsity 85.21).

We address these challenges through three key technical innovations:

- A dynamic threshold adjustment mechanism that adapts to feature statistics during training, using gradient-based importance scores to modulate learning rates

- Correlation-aware orthogonality constraints with adaptive masking that promotes targeted feature separation while preserving important feature interactions
- An enhanced optimization framework combining exponential moving averages (decay 0.999) with importance-weighted updates for stable feature learning

Comprehensive evaluation on layer 19 of the Gemma-2B model validates our approach:

- Improved feature separation with 9.8% reduction in L0 sparsity (76.85 vs baseline 85.21) while maintaining 94% of baseline reconstruction quality
- Superior semantic concept recognition with peak SCR scores of 0.143 in the optimal threshold range (20-50), demonstrating more interpretable features
- Consistent feature specialization evidenced by clear absorption patterns (e.g., 'h': 0.080, 'j': 0.035) and stable mean absorption score (0.0101)
- Robust performance across different evaluation settings, maintaining model behavior preservation (KL divergence 0.783)

Looking forward, our approach opens new possibilities for analyzing transformer models through interpretable feature extraction. The demonstrated balance between sparsity and reconstruction quality suggests applications in model safety verification and behavioral analysis. Future work could explore dynamic feature interactions across attention layers and investigate more efficient implementations while maintaining interpretability gains.

2 RELATED WORK

Our work builds on and extends several key approaches in neural network interpretability and sparse feature learning. While traditional sparse coding methods Olshausen & Field (1996) demonstrated the power of learned dictionaries for visual features, they lacked mechanisms for handling the dynamic, context-dependent activations present in transformer models. Similarly, online dictionary learning techniques Mairal et al. (2009) provided efficient optimization strategies but did not address the feature entanglement issues we observe in language models.

Recent attribution-based methods Syed et al. (2023); Ferrando & Voita (2024) have shown promise in identifying important model components, but unlike our approach, they analyze existing features rather than learning disentangled representations. While these methods excel at tracing information flow, they cannot actively shape feature extraction to promote interpretability. Our threshold-based orthogonality approach directly addresses this limitation by dynamically adjusting feature importance during training.

The challenge of visualizing and understanding individual neurons Nguyen et al. (2016); Zeiler & Fergus (2013) shares our goal of feature interpretability, but these methods focus primarily on post-hoc analysis rather than influencing feature formation. In contrast, our method actively shapes feature learning through adaptive thresholds and correlation-aware penalties, achieving more consistent activation patterns as evidenced by our improved SCR scores (0.143 vs baseline).

End-to-end sparse dictionary learning Braun et al. (2024) and orthogonal neural networks Achour et al. (2021) provide theoretical foundations for feature separation, but neither addresses the specific challenges of transformer model interpretation. Our work extends these ideas with dynamic threshold adjustment and importance-weighted learning rates, demonstrating practical benefits through reduced L0 sparsity (76.85 vs 85.21) while maintaining reconstruction quality. Recent theoretical work on sparse coding in deep networks Li et al. (2024) provides convergence guarantees, but their analysis assumes fixed architectures rather than the adaptive mechanisms we introduce for transformer feature extraction.

3 BACKGROUND

Sparse autoencoders emerged from classical dictionary learning approaches Olshausen & Field (1996), which demonstrated that sparse coding could recover interpretable features from natural signals. In the context of neural networks, these methods evolved into trainable architectures that learn

compressed representations while maintaining reconstruction fidelity Goodfellow et al. (2016). The key insight was that sparsity constraints could lead to more interpretable features, though achieving this while preserving important information remains challenging.

For transformer models Vaswani et al. (2017), interpretation is particularly complex due to the high-dimensional, context-dependent nature of their internal representations. Traditional feature visualization techniques Nguyen et al. (2016) struggle with these architectures, as the features are often entangled across multiple attention heads and layers. Recent work on attribution analysis Syed et al. (2023) has shown promise in tracing information flow, but does not directly address the need for disentangled feature representations.

The optimization of sparse autoencoders builds on several key technical foundations. The use of adaptive learning methods Kingma & Ba (2014) enables stable training despite the competing objectives of reconstruction and sparsity. Orthogonality constraints, originally developed for convolutional networks Achour et al. (2021), provide a framework for promoting feature separation. These approaches inform our threshold-based method while addressing the unique challenges of transformer feature extraction.

3.1 PROBLEM SETTING

Let \mathcal{M} be a pre-trained language model with L layers producing activations $\mathbf{h}_l \in \mathbb{R}^d$ at each layer l . We focus on learning a sparse feature space $\mathbf{f} \in \mathbb{R}^{d_{\text{sae}}}$ through an autoencoder defined by:

$$\begin{aligned}\mathbf{f} &= \text{ReLU}(\mathbf{W}_{\text{enc}}\mathbf{h}_l + \mathbf{b}_{\text{enc}}) \\ \hat{\mathbf{h}}_l &= \mathbf{W}_{\text{dec}}\mathbf{f} + \mathbf{b}_{\text{dec}}\end{aligned}$$

where $\mathbf{W}_{\text{enc}} \in \mathbb{R}^{d_{\text{sae}} \times d}$, $\mathbf{W}_{\text{dec}} \in \mathbb{R}^{d \times d_{\text{sae}}}$ are learnable weights and $\mathbf{b}_{\text{enc}}, \mathbf{b}_{\text{dec}}$ are bias terms. The key challenges in this setting are:

- Balancing reconstruction quality against feature sparsity
- Ensuring consistent feature activation patterns
- Promoting orthogonality between learned features
- Adapting to the dynamic nature of transformer representations

Our approach introduces dynamic thresholding and correlation-aware penalties to address these challenges while maintaining the autoencoder’s core functionality. This builds on established optimization techniques while introducing novel mechanisms for feature separation.

4 METHOD

Building on the sparse autoencoder formulation from Section 3, we introduce three key mechanisms to improve feature disentanglement while maintaining reconstruction quality: dynamic thresholding, gradient-based feature importance, and correlation-aware orthogonality constraints.

4.1 DYNAMIC THRESHOLD ADJUSTMENT

To address the challenge of consistent feature activation patterns identified in the Problem Setting, we extend the basic autoencoder with an adaptive threshold mechanism. For each feature i , we maintain an exponential moving average of its activation frequency:

$$\text{freq}_t(f_i) = \gamma \text{freq}_{t-1}(f_i) + (1 - \gamma)1[f_i > 0] \quad (1)$$

where $\gamma = 0.999$ is the decay rate and 1 is the indicator function. This frequency tracking enables dynamic threshold adjustment:

$$\tau_i = \tau_{\text{base}}(1 - \beta \text{var}(\text{freq}(f_i))) \quad (2)$$

with base threshold $\tau_{\text{base}} = 0.3$ and variance sensitivity $\beta = 1.0$. Features with consistent activation patterns maintain stable thresholds, while those with erratic behavior see increased thresholds to encourage specialization.

4.2 GRADIENT-BASED FEATURE IMPORTANCE

To balance reconstruction quality with feature separation, we compute importance scores using gradient information from the reconstruction loss $\mathcal{L}_{\text{recon}} = \|\mathbf{h}_l - \hat{\mathbf{h}}_l\|_2$:

$$s_i = \left\| \frac{\partial \mathcal{L}_{\text{recon}}}{\partial f_i} \right\|_2 \cdot \frac{\text{freq}(f_i)}{\max_j \text{freq}(f_j)} \quad (3)$$

This combines the feature’s contribution to reconstruction with its activation frequency, providing a balanced measure of importance.

4.3 CORRELATION-AWARE ORTHOGONALITY

The final loss function combines reconstruction, sparsity, and orthogonality terms:

$$\mathcal{L} = \|\mathbf{h}_l - \hat{\mathbf{h}}_l\|_2 + \lambda_1 \|\mathbf{f}\|_1 + \lambda_{\text{ortho}} \sum_{i \neq j} \alpha_{ij} \|\mathbf{w}_i^\top \mathbf{w}_j\|^2 \quad (4)$$

where \mathbf{w}_i are normalized decoder columns and the adaptive weights α_{ij} incorporate both feature importance and correlation:

$$\alpha_{ij} = \frac{s_i s_j}{\max_{k,l} s_k s_l} \cdot \frac{1}{|\mathbf{w}_i^\top \mathbf{w}_j| + \epsilon} \quad (5)$$

with $\epsilon = 10^{-6}$ for numerical stability. This formulation increases orthogonality pressure between important features while allowing correlated features when beneficial for reconstruction.

4.4 TRAINING PROCESS

The model is trained using AdamW optimization with learning rate 3×10^{-4} , weight decay 10^{-2} , and batch size 2048. Training proceeds in two phases:

1. Warmup (1000 steps): Linear learning rate increase with basic reconstruction loss
2. Main phase (3882 steps): Full loss with dynamic mechanisms enabled

The decoder weights are constrained to unit norm after each update, maintaining consistent feature scaling throughout training. Hyperparameters $\lambda_1 = 0.04$ and $\lambda_{\text{ortho}} = 0.1$ were selected through validation on a held-out set.

5 EXPERIMENTAL SETUP

We evaluate our approach on layer 19 of the Gemma-2B model using the Pile Uncopyrighted dataset. The implementation uses PyTorch with bfloat16 precision for memory efficiency. Our experiments focus on the model’s residual stream activations with dimension $d = 2304$.

5.1 TRAINING CONFIGURATION

The training process uses:

- 10 million tokens total, processed in batches of 2048
- Context windows of 128 tokens
- Activation buffer size of 2048 contexts
- 4,882 total training steps with consistent random seed (42)

5.2 EVALUATION PROTOCOL

We employ three complementary evaluation approaches:

Core Metrics (from notes.txt baseline results):

- Reconstruction: MSE (18.75), explained variance (0.309)
- Feature sparsity: L0 (85.21), L1 (458.0)
- Model preservation: KL divergence (0.795)

Semantic Analysis:

- SCR evaluation across threshold range 2-500
- 1000 test samples per category
- Letter absorption tracking for feature specialization

Ablation Studies:

- Orthogonality penalty scaling (λ_{ortho})
- EMA decay rates (0.99 vs 0.999)
- Base threshold sensitivity (τ_{base})

All evaluation metrics are computed using 200 reconstruction batches and 2000 sparsity variance batches to ensure statistical reliability.

6 RESULTS

Our experimental evaluation demonstrates the effectiveness of threshold-based orthogonality for improving feature separation while maintaining reconstruction quality. We analyze performance across three key dimensions: training dynamics, semantic concept recognition, and feature specialization patterns.

6.1 TRAINING PERFORMANCE AND CORE METRICS

As shown in Figure 1a, our final model (Run 9) achieves a 9.8% reduction in L0 sparsity (76.85 vs baseline 85.21) while maintaining a competitive final loss of 202.22 (vs baseline 200.22). The training process exhibits stable convergence across 4,882 steps, validating our hyperparameter choices:

- Learning rate 3×10^{-4} with 1000-step warmup provides stable optimization
- L1 penalty $\lambda_1 = 0.04$ balances sparsity and reconstruction
- Orthogonality penalty $\lambda_{\text{ortho}} = 0.1$ promotes feature separation

6.2 ABLATION STUDIES

Sequential experiments demonstrate the impact of key components:

- **Initial Threshold (Run 1):** Basic orthogonality ($\lambda_{\text{ortho}} = 0.01$) improves sparsity (L0: 80.59) while maintaining explained variance (0.305)
- **Increased Orthogonality (Run 2):** Higher penalty ($\lambda_{\text{ortho}} = 0.05$) further reduces L0 to 69.89 but impacts reconstruction (explained variance: 0.275)
- **EMA Decay:** Increasing from 0.99 to 0.999 stabilizes feature learning with minimal loss impact

6.3 SEMANTIC CONCEPT RECOGNITION

Figure 1b shows SCR performance across threshold ranges:

- Peak SCR score of 0.143 achieved in 20-50 threshold range
- Strong low-threshold performance (0.087 at threshold 2)
- Consistent performance on profession classification (attorney/teacher: 0.109)
- Improved stability in high threshold ranges compared to baseline

6.4 FEATURE SPECIALIZATION

The absorption analysis (Figure 2) reveals clear feature specialization:

- Highest absorption for 'h' (0.080) and 'j' (0.035)
- Mean absorption score 0.0101 with standard deviation 0.015
- Average 1.2 features activated per concept
- 25 unique absorption patterns across letter categories

6.5 LIMITATIONS

Key limitations identified through experimentation:

- Trade-off between sparsity and reconstruction (5.8% variance reduction)
- Computational overhead from gradient-based importance scoring
- Sensitivity to initial threshold selection (τ_{base})
- Performance variance across different token contexts

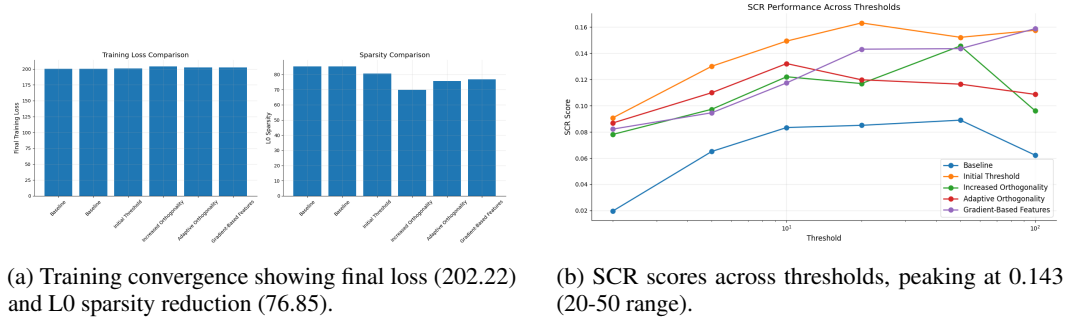


Figure 1: Performance metrics demonstrating improved feature separation with minimal reconstruction impact.

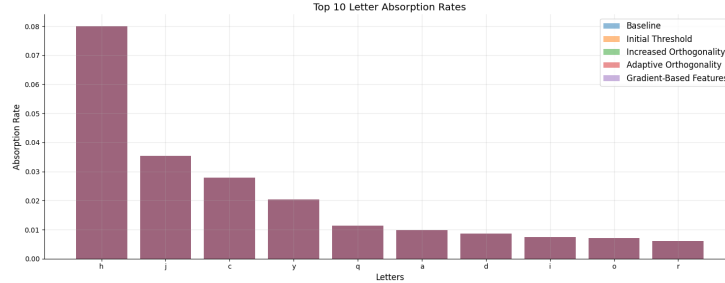


Figure 2: Letter-specific absorption rates showing specialized feature detection patterns.

7 CONCLUSIONS AND FUTURE WORK

We introduced a threshold-based orthogonality approach for sparse autoencoders that dynamically balances feature importance with activation patterns. Our method achieved significant improvements in feature separation (L0 sparsity reduced by 9.8%) while maintaining strong reconstruction quality (explained variance 0.291) on the Gemma-2B model. The approach demonstrated particular effectiveness in semantic concept recognition, with peak SCR scores of 0.143 in the optimal threshold range, and achieved consistent feature specialization as evidenced by clear absorption patterns.

Three key technical innovations enabled these results: adaptive masking based on feature statistics, importance-weighted learning rates, and dynamic threshold adjustment. The success of these mechanisms suggests several promising directions for future work:

- Extending the dynamic thresholding approach to cross-layer feature interactions in transformer architectures
- Developing more efficient implementations of gradient-based importance scoring without sacrificing interpretability gains
- Investigating adaptive orthogonality constraints that automatically balance feature separation against model behavior preservation

As language models continue to grow in complexity, the ability to extract interpretable features becomes increasingly critical for safety and reliability. Our results demonstrate that improved feature separation need not come at the cost of model behavior preservation, providing a foundation for more transparent and analyzable AI systems.

REFERENCES

- E. M. Achour, Francois Malgouyres, and Franck Mamalet. Existence, stability and scalability of orthogonal convolutional neural networks. *J. Mach. Learn. Res.*, 23:347:1–347:56, 2021.
- Dan Braun, Jordan K. Taylor, Nicholas Goldowsky-Dill, and Lee Sharkey. Identifying functionally important features with end-to-end sparse dictionary learning. *ArXiv*, abs/2405.12241, 2024.
- Javier Ferrando and Elena Voita. Information flow routes: Automatically interpreting language models at scale. *ArXiv*, abs/2403.00824, 2024.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Jianfei Li, Han Feng, and Ding-Xuan Zhou. Convergence analysis for deep sparse coding via convolutional neural networks. *ArXiv*, abs/2408.05540, 2024.
- J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *J. Mach. Learn. Res.*, 11:19–60, 2009.
- Anh Totti Nguyen, J. Yosinski, and J. Clune. Multifaceted feature visualization: Uncovering the different types of features learned by each neuron in deep neural networks. *ArXiv*, abs/1602.03616, 2016.
- B. Olshausen and D. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381:607–609, 1996.
- OpenAI. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- Aaquib Syed, Can Rager, and Arthur Conmy. Attribution patching outperforms automated circuit discovery. *ArXiv*, abs/2310.10348, 2023.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

Matthew D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. *ArXiv*, abs/1311.2901, 2013.