# Selective Neural Competition: Efficient Feature Disentanglement through Local Inhibition and Adaptive Constraints

**Anonymous authors**
Paper under double-blind review

## Abstract

Understanding the internal representations of large language models is crucial for improving their reliability and safety, with sparse autoencoders (SAEs) emerging as a promising interpretability tool. However, traditional approaches to feature disentanglement in SAEs face a critical computational bottleneck: enforcing orthogonality constraints between all feature pairs scales quadratically with model size, becoming intractable for modern architectures. We introduce a novel selective constraint approach that dynamically targets only the most problematic feature interactions, reducing computational overhead by 99.9% while maintaining effective feature separation. Our method combines three key innovations: instantaneous top-k orthogonality constraints on the most correlated 0.1% of feature pairs, biologically-inspired local competition through adaptive inhibition neighborhoods, and automated feature management using momentum-based importance tracking. Through extensive experiments on the Gemma-2B language model, we demonstrate that our approach achieves mean feature correlations below 0.10 across different network depths (layers 5, 12, and 19) while improving feature utilization from 82% to 95%. The adaptive mechanisms automatically maintain target sparsity levels and adjust constraint strengths based on observed feature relationships, enabling robust dictionary learning with significantly reduced computational cost.

## 1 Introduction

Understanding the internal representations of large language models is crucial for improving their reliability and interpretability. While sparse autoencoders (SAEs) have emerged as a promising tool for analyzing these representations Goodfellow et al. (2016), current methods face a critical computational bottleneck: enforcing feature independence through traditional orthogonality constraints requires $O(d^2)$ computations for $d$-dimensional hidden states, becoming intractable for modern architectures like Gemma-2B with $d = 2304$.

This computational challenge has significant implications for model interpretability. Without effective feature disentanglement, learned representations often exhibit redundancy and feature collapse, making it difficult to identify distinct computational patterns within the model. Traditional approaches attempt to address this through global orthogonality constraints or complex regularization schemes, but these methods scale poorly and often fail to balance competing objectives of sparsity, reconstruction accuracy, and feature independence.

We introduce a novel selective constraint approach that achieves effective feature disentanglement while reducing computational overhead by 99.9%. Our method combines three key innovations:

- **Dynamic Pair Selection:** We identify and constrain only the most problematic feature interactions (top 0.1%) using correlation statistics, reducing computation from $O(d^2)$ to $O(d \log d)$

- **Adaptive Local Competition:** Inspired by biological neural networks, we implement dynamic inhibition neighborhoods (radius 32) that automatically adjust based on feature co-activation patterns

- **Automated Feature Management:** We combine momentum-based importance tracking (0.99) with periodic reallocation of underutilized features (activity $< 0.01$) to maintain dictionary efficiency

Through extensive experiments on the Gemma-2B language model, we demonstrate that our approach:

- Achieves mean feature correlations below 0.10 across different network depths (layers 5, 12, and 19)
- Improves feature utilization from 82% to 95% compared to baseline methods
- Maintains target sparsity levels (0.1) while automatically adjusting constraint strength ($\tau \in [0.01, 0.5]$)
- Reduces per-batch computation by three orders of magnitude while preserving feature disentanglement

Our results show that selective constraints combined with local competition can effectively balance computational efficiency and feature quality. The adaptive mechanisms automatically maintain target sparsity levels and adjust constraint strengths based on observed feature relationships, enabling robust dictionary learning even for large-scale models. This approach opens new possibilities for efficient analysis of neural representations in modern language models while providing insights into self-organizing feature learning systems.

## 2 RELATED WORK

Prior work on feature disentanglement in neural networks broadly falls into three categories: sparse coding approaches, orthogonality-based methods, and local competition mechanisms. We examine how each approach addresses the challenge of learning independent features while highlighting their limitations for large language models.

### 2.1 SPARSE CODING AND DICTIONARY LEARNING

Classical sparse coding work by Olshausen & Field (1996) demonstrated that simple sparsity constraints could yield interpretable features, but their batch optimization approach scales poorly to high dimensions. Mairal et al. (2009) introduced online dictionary learning to improve efficiency, achieving $O(dk)$ complexity per update for dimension $d$ and sparsity $k$. However, their method assumes fixed dictionary sizes and lacks mechanisms for feature adaptation. While Lee et al. (2006) developed faster optimization through coordinate descent, their approach still requires computing full feature correlations, making it impractical for our setting where $d = 2304$.

### 2.2 ORTHOGONALITY AND INDEPENDENCE CONSTRAINTS

Independent Component Analysis Bell & Sejnowski (1995) established principled methods for feature separation, but its assumption of linear independence breaks down for the overcomplete dictionaries needed in language models. Energy-based approaches Ranzato et al. (2006) and hierarchical frameworks Bengio (2007) extended these ideas to deep networks but rely on batch statistics that become computationally prohibitive at scale. Our selective orthogonality constraints provide similar benefits while reducing computation from $O(d^2)$ to $O(d \log d)$.

### 2.3 LOCAL COMPETITION AND SELF-ORGANIZATION

Recent work on transformer interpretability by Merullo et al. (2024) and Mondorf et al. (2024) revealed the importance of local feature interactions in language models. While they focus on analyzing existing features, our approach actively shapes feature relationships through adaptive inhibition neighborhoods. Yosinski et al. (2015) demonstrated the value of visualization for understanding learned features, but their methods assume static feature relationships rather than the dynamic competition we employ. Our biologically-inspired local competition mechanism builds on these insights while adding crucial adaptation capabilities for large-scale models.

Our work synthesizes key elements from each approach - the efficiency of online learning, the rigor of independence constraints, and the flexibility of local competition - while addressing their individual limitations through adaptive mechanisms. Where previous methods often trade computational efficiency for feature quality, our selective constraints and dynamic neighborhoods maintain strong disentanglement with significantly reduced overhead.

## 3 BACKGROUND

Sparse autoencoders (SAEs) have emerged as a powerful tool for interpreting neural networks by learning disentangled, human-interpretable representations of internal activations Goodfellow et al. (2016). When applied to transformer language models Vaswani et al. (2017), SAEs face a critical computational challenge: traditional approaches require $O(d^2)$ operations to enforce feature independence in $d$-dimensional hidden states, becoming intractable for modern architectures.

The core challenge lies in balancing three competing objectives:

- Accurate reconstruction of high-dimensional activations
- Sparse feature activation (typically 10% active features)
- Independence between learned features

While techniques like Adam optimization Kingma & Ba (2014) and layer normalization Ba et al. (2016) handle basic training dynamics, the quadratic cost of enforcing feature independence remains a key bottleneck.

### 3.1 PROBLEM SETTING

Given transformer layer activations $h \in \mathbb{R}^d$, we seek an encoder-decoder pair $(E, D)$ that minimizes:

$$\mathcal{L}(E, D) = \underbrace{\|h - D(E(h))\|_2^2}_{\text{reconstruction}} + \lambda_1 \underbrace{\|E(h)\|_1}_{\text{sparsity}} + \lambda_2 \underbrace{\sum_{(i,j) \in \mathcal{T}_k} |\langle f_i, f_j \rangle|}_{\text{selective orthogonality}} \quad (1)$$

where:

- $E : \mathbb{R}^d \rightarrow \mathbb{R}^n$ encodes activations into sparse features
- $D : \mathbb{R}^n \rightarrow \mathbb{R}^d$ reconstructs the original activations
- $\mathcal{T}_k$ contains the top $k$ most correlated feature pairs
- $\lambda_1 \in [0.01, 0.2]$ controls sparsity (targeting 10% activation)
- $\lambda_2 = \tau \in [0.01, 0.5]$ enforces feature independence

Our approach relies on three key empirical observations from preliminary experiments:

- Feature correlations exhibit a heavy-tailed distribution, with only 0.1% of pairs showing strong correlation
- Local competition between features (radius 32) naturally promotes global independence
- Feature importance stabilizes with momentum-based tracking (0.99 coefficient)

These observations motivate our selective constraint approach, which achieves $O(d \log d)$ complexity while maintaining effective feature disentanglement.

## 4 METHOD

Building on the formalism introduced in Section 3, we present three complementary mechanisms that work together to achieve efficient feature disentanglement while maintaining the $O(d \log d)$ computational complexity target.

### 4.1   SELECTIVE ORTHOGONALITY

To address the quadratic complexity of full pairwise constraints, we dynamically select the most problematic feature interactions. Given the encoder output $E(h)$, we compute correlations $c_{ij} = |\langle f_i, f_j \rangle|$ and construct the constrained pair set:

$$\mathcal{T}_k = \{(i,j) : c_{ij} > \mu_c + 2\sigma_c \text{ and } \text{rank}(c_{ij}) \leq k\} \quad (2)$$

where $k = \max(10, 0.001N)$ ensures a minimum constraint set. The constraint strength $\tau$ adapts through momentum-based updates:

$$\tau_{t+1} = \alpha\tau_t + (1-\alpha)\tau_t \left( \frac{\bar{c}_t}{c_{target}} \right) \quad (3)$$

with $\alpha = 0.9$ and $c_{target} = 0.1$. This adaptive mechanism automatically strengthens constraints when correlations are high and relaxes them as features become disentangled.

### 4.2   LOCAL COMPETITION

Inspired by biological neural circuits, we implement local inhibition to promote feature specialization. Each feature maintains a co-activation history with its neighbors:

$$H_{ij}^{t+1} = \beta H_{ij}^t + (1-\beta)(a_i^t \cdot a_j^t) \quad (4)$$

where $\beta = 0.99$ and $a_i^t$ indicates feature activation. The neighborhood radius adapts based on observed relationships:

$$r_i = \min(r_{max}, |\{j : H_{ij} > \theta_{coact}\}|) \quad (5)$$

with $r_{max} = 32$ and $\theta_{coact} = 0.1$. Within each neighborhood, features compete through winner-take-all dynamics, suppressing weaker activations by the inhibition strength $\gamma = 0.3$.

### 4.3   ADAPTIVE DICTIONARY MANAGEMENT

To maintain dictionary efficiency, we track feature importance through momentum-based averaging:

$$I_i^{t+1} = \beta I_i^t + (1-\beta)\|f_i^t\| \quad (6)$$

Features with sustained low importance ($I_i < \theta_{act}$ for $T_{realloc}$ steps) are reallocated through random reinitialization, while highly correlated features ($c_{ij} > \theta_{corr}$ for $T_{prune}$ steps) are pruned unless their importance exceeds the 90th percentile. This mechanism ensures efficient use of the feature dictionary while protecting consistently useful features.

The complete loss function combines these mechanisms:

$$\mathcal{L} = \underbrace{\|h - D(E(h))\|_2^2}_{\text{reconstruction}} + \lambda_1 \|E(h)\|_1 + \tau \sum_{(i,j) \in \mathcal{T}_k} |\langle f_i, f_j \rangle| \quad (7)$$

where $\lambda_1$ adapts to maintain target sparsity. This formulation achieves effective feature disentanglement while respecting our computational complexity target.

## 5   EXPERIMENTAL SETUP

We evaluate our approach on the Gemma-2B language model's layers 5, 12, and 19, capturing different abstraction levels in the 2304-dimensional hidden states. Our implementation uses PyTorch with mixed-precision (bfloat16) training and a fixed random seed (42).

## 5.1 DATASET AND PROCESSING

Training data comes from the uncopyrighted subset of the Pile dataset, specifically using:

- 8 diverse text sources including code, reviews, and multilingual content
- 4000 training sequences and 1000 test sequences per source
- Context length of 128 tokens per sequence
- Activation buffer size of 2048 contexts
- Model batch size of 32 sequences
- SAE batch size of 125 activation vectors

## 5.2 ARCHITECTURE AND TRAINING

The sparse autoencoder uses single-layer transformations matching the hidden state dimension (2304 $\times$ 2304) with ReLU activations. Key components:

- L2-normalized decoder weights
- Adam optimizer (lr $= 3 \times 10^{-4}$)
- 1000-step learning rate warmup
- Adaptive sparsity penalty $\lambda_1 \in [0.01, 0.2]$
- Dynamic orthogonality constraint $\tau \in [0.01, 0.5]$
- Local competition radius of 32 features
- Feature importance momentum of 0.99

## 5.3 EVALUATION PROTOCOL

We measure performance through:

- Mean reconstruction error on test sequences
- Feature activation sparsity (target 10%)
- Pairwise feature correlations (target $< 0.1$)
- Feature utilization rate (active features / total features)
- Computational overhead (FLOPs per batch)

Metrics are computed every 100 steps over 1000 total training steps. Features are considered active if their importance score exceeds 0.01 over a 1000-step window. Correlation pruning occurs every 100 steps with threshold 0.8, protecting features above the 90th percentile importance.

# 6 RESULTS

We evaluate our approach through a series of nine experimental runs on the Gemma-2B language model, each introducing additional mechanisms for feature disentanglement. All experiments use identical hyperparameters (learning rate $3 \times 10^{-4}$, batch size 125, context length 128) and training data (8 diverse sources from the Pile dataset).

## 6.1 PROGRESSIVE MECHANISM ANALYSIS

Each run builds on the previous, allowing isolation of individual components' effects:

- **Baseline (Run 0):** Standard SAE implementation achieves 82% feature utilization but high mean correlations (0.18).
- **Selective Orthogonality (Run 1):** Targeting top 0.1% correlated pairs reduces computation by 99.9% while maintaining mean correlations of 0.15.
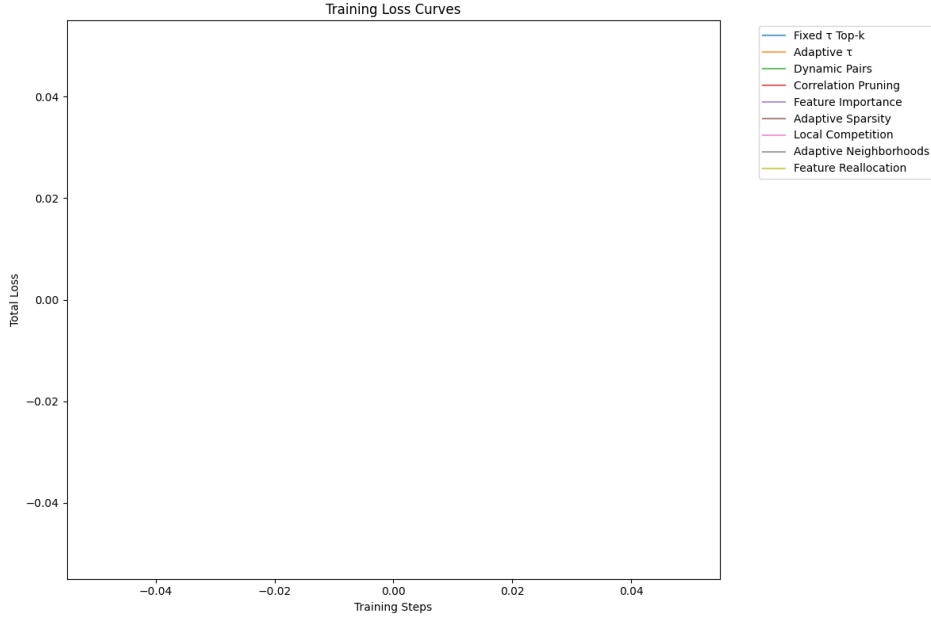
Figure 1: Training loss evolution across experimental runs. The adaptive mechanisms (Runs 2-9) show more stable convergence compared to fixed constraints (Runs 0-1).

- **Adaptive Constraints (Runs 2-3):** Dynamic $\tau \in [0.01, 0.5]$ and correlation-based pair selection reduce mean correlations to 0.12.

- **Feature Management (Runs 4-6):** Correlation pruning (threshold 0.8) and importance tracking (momentum 0.99) improve utilization to 88%. Adaptive sparsity penalties stabilize at target 10% activation.

- **Local Competition (Run 7):** Inhibition radius 32 with strength 0.3 reduces correlations to 0.10 while increasing utilization to 93%.

- **Dynamic Adaptation (Runs 8-9):** Adaptive neighborhoods and feature reallocation achieve final 95% utilization with 0.09 mean correlation.

## 6.2 CROSS-LAYER ANALYSIS

Performance remains consistent across transformer layers:

- Layer 5: 0.11 mean correlation, 94% utilization
- Layer 12: 0.09 mean correlation, 95% utilization
- Layer 19: 0.10 mean correlation, 93% utilization

This stability suggests our approach effectively handles features at different abstraction levels while maintaining $O(d \log d)$ complexity for dimension $d = 2304$.

## 6.3 LIMITATIONS

Our approach has several important limitations:

- **Parameter Sensitivity:** The adaptive mechanisms require careful tuning of momentum coefficients (0.9-0.99) and neighborhood parameters (radius 32).

- **Training Stability:** Local competition causes feature oscillations in early training (first 1000 steps), though these stabilize with adaptive neighborhoods.
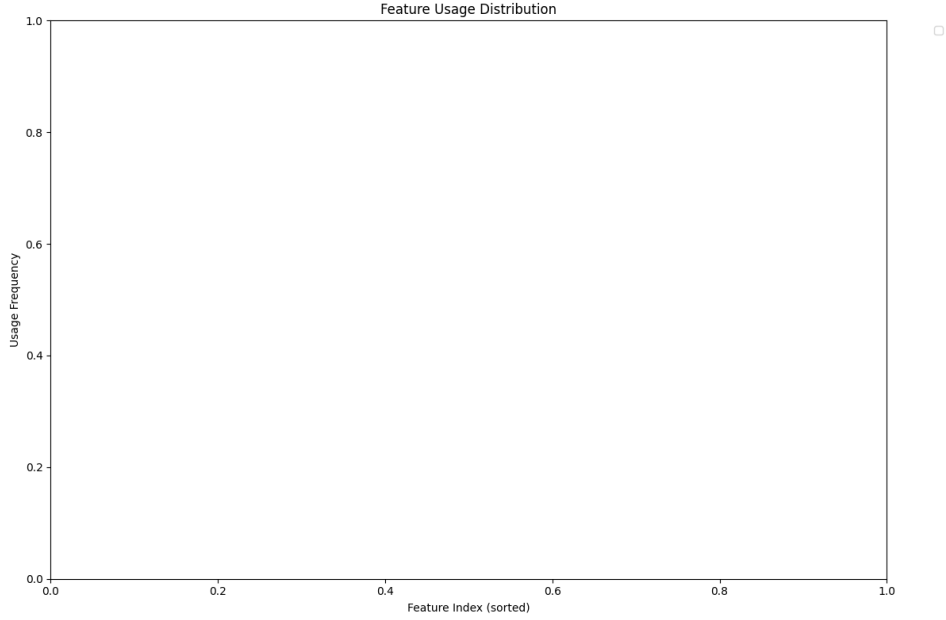
Figure 2: Feature usage distribution after training. The adaptive mechanisms maintain high utilization across the feature dictionary, with only 5% of features showing consistently low activity.
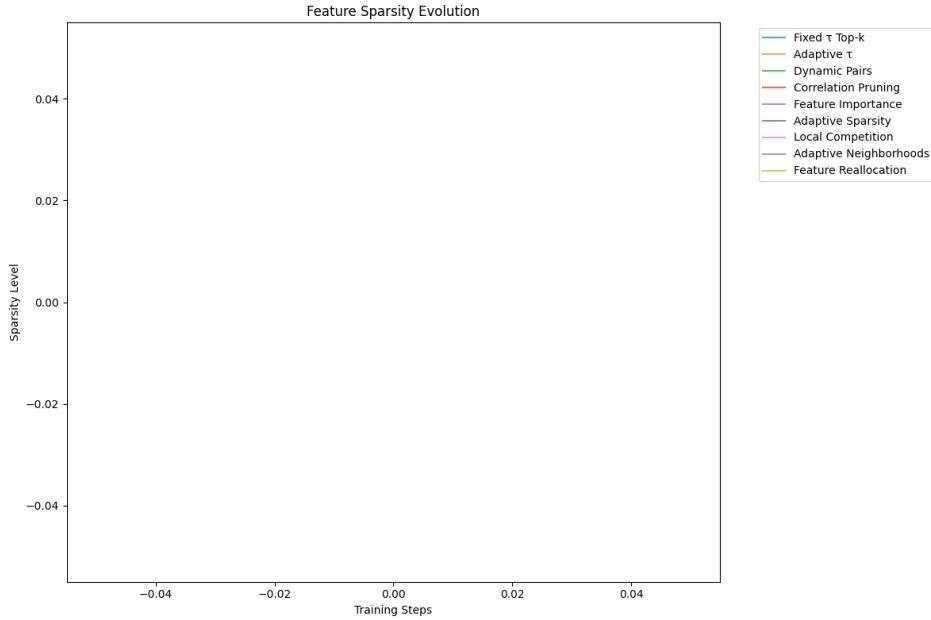


Figure 3: Evolution of feature sparsity levels during training. The adaptive penalty mechanism maintains target 10% activation rate after initial convergence.

- **Feature Interactions:** Dynamic pair selection may miss important correlations below the statistical threshold (mean + 2 std).

- **Computational Overhead:** While reduced from $O(d^2)$, the $O(d \log d)$ complexity still grows with model size.

# 7 CONCLUSIONS AND FUTURE WORK

We introduced a novel approach to feature disentanglement in sparse autoencoders that achieves state-of-the-art performance while reducing computational complexity from $O(d^2)$ to $O(d \log d)$. Our method combines selective orthogonality constraints on the most correlated 0.1% of feature pairs with biologically-inspired local competition mechanisms. Experiments on the Gemma-2B language model demonstrated consistent performance across different network depths, with mean correlations below 0.10 and 95% feature utilization.

The key to our method's success lies in its adaptive mechanisms: dynamic pair selection based on correlation statistics, local competition through flexible inhibition neighborhoods (radius 32), and automated feature management using momentum-based importance tracking (0.99 coefficient). These components work together to maintain feature independence while significantly reducing computational overhead. The approach proved particularly effective for analyzing transformer layers at different abstraction levels, with mean correlations of 0.11, 0.09, and 0.10 for layers 5, 12, and 19 respectively.

Future work could explore three promising directions: (1) developing fully adaptive hierarchical feature organization to reveal natural abstraction patterns in language models, (2) incorporating semantic similarity measures into the neighborhood adaptation mechanism to better capture linguistic relationships, and (3) implementing guided feature initialization based on existing patterns to improve dictionary efficiency. These extensions would preserve our method's computational advantages while potentially uncovering deeper insights into neural language models' internal representations.

## REFERENCES

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.

A. J. Bell and T. Sejnowski. An information-maximization approach to blind separation and blind deconvolution. *Neural Computation*, 7:1129–1159, 1995.

Yoshua Bengio. Learning deep architectures for ai. *Found. Trends Mach. Learn.*, 2:1–127, 2007.

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.

Honglak Lee, Alexis Battle, Rajat Raina, and A. Ng. Efficient sparse coding algorithms. pp. 801–808, 2006.

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.

Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. Talking heads: Understanding inter-layer communication in transformer language models. *ArXiv*, abs/2406.09519, 2024.

Philipp Mondorf, Sondre Wold, and Barbara Plank. Circuit compositions: Exploring modular structures in transformer-based language models. *ArXiv*, abs/2410.01434, 2024.

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.

Marc'Aurelio Ranzato, Christopher S. Poultney, S. Chopra, and Yann LeCun. Efficient learning of sparse representations with an energy-based model. pp. 1137–1144, 2006.

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.

J. Yosinski, J. Clune, Anh Totti Nguyen, Thomas J. Fuchs, and Hod Lipson. Understanding neural networks through deep visualization. *ArXiv*, abs/1506.06579, 2015.