# Competitive Progressive Autoencoders: Learning Sparse Features Through Adaptive Activation

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

## Abstract

Interpreting large language models through sparse autoencoders (SAEs) is hindered by feature collapse and inefficient capacity utilization, where features either become redundant or remain underutilized during training. We present a progressive feature learning framework that addresses these challenges through two key innovations: adaptive feature enabling that starts with 25% capacity and incrementally activates features based on utilization metrics, and temperature-scaled competition (range 0.05-0.3) that prevents feature collapse through double exponential density scaling. Experiments on the Gemma-2B model demonstrate a 53.6% improvement in sparsity (L0: 42.05 vs baseline 90.60) while maintaining strong reconstruction quality (MSE: 28.63). Our hybrid increment strategy achieves 12.4% better feature utilization compared to fixed increments, with particularly strong performance in selective concept response at low thresholds (SCR: 0.027). The architecture's effectiveness is further validated through stable model preservation scores (KL divergence: 0.481) and consistent feature absorption metrics (0.0101), demonstrating robust feature specialization without sacrificing interpretability.

## 1 Introduction

The challenge of interpreting large language models (LLMs) has become increasingly critical as these models grow in complexity and capability OpenAI (2024). While sparse autoencoders (SAEs) offer a promising approach for extracting interpretable features from neural activations Goodfellow et al. (2016), current methods struggle with two fundamental challenges: feature collapse, where multiple dictionary elements encode redundant patterns, and inefficient capacity utilization, where features remain inactive throughout training.

These challenges arise from the inherent tension between sparsity and representational power. Traditional approaches either activate all features simultaneously, leading to interference and redundancy, or employ fixed sparsity constraints that can prevent the capture of important patterns. The dynamic nature of feature importance during training further complicates this balance - features that are crucial early in training may become redundant later, while initially quiet features may eventually encode essential patterns.

We address these challenges through a novel progressive feature learning framework that combines two key mechanisms:

- **Adaptive Feature Enabling:** Starting with 25% of total capacity, our approach incrementally activates features using a hybrid strategy based on utilization metrics. This includes fixed increments (5 features) at low utilization (<30%), adaptive scaling (2-7 features) at medium utilization (30-60%), and aggressive growth (4-10 features) at high utilization (60%).

- **Temperature-Scaled Competition:** A momentum-based temperature scaling mechanism (range 0.05-0.3) prevents feature collapse through double exponential density scaling, allowing features to specialize while maintaining stability.

Experiments on the Gemma-2B model demonstrate significant improvements across key metrics:

- **Sparsity:** 53.6% improvement in L0 sparsity (42.05 vs baseline 90.60) and 64.2% reduction in L1 penalty (168.75 vs 471.54)
- **Feature Specialization:** Strong selective concept response (SCR) at low thresholds (0.027) while maintaining feature absorption metrics (0.0101)
- **Model Preservation:** Stable KL divergence scores (0.481) despite increased sparsity

Our analysis reveals that starting with fewer features but allowing faster growth achieves better feature selectivity, as evidenced by improved SCR metrics across all thresholds. The hybrid increment strategy combined with temperature annealing proves crucial for balancing feature competition and preservation, though some trade-off in reconstruction quality (MSE: 28.63 vs baseline 18.63) is observed.

Looking ahead, three promising directions emerge: (1) semantic-aware feature enabling that considers conceptual relationships, (2) hierarchical competition mechanisms operating at multiple granularities, and (3) extension to other architectures beyond language models. Our results suggest that progressive feature learning could become a fundamental approach for developing interpretable representations across deep learning domains.

## 2 BACKGROUND

The foundations of sparse feature learning trace back to neuroscience-inspired work by Olshausen & Field (1996), who demonstrated how simple cells in the visual cortex learn efficient coding through competitive mechanisms. This was extended to artificial neural networks through dictionary learning approaches Kreutz-Delgado et al. (2003), which established the importance of overcomplete representations for capturing complex patterns.

In the context of modern language models, interpretability has become increasingly critical as model complexity grows Vaswani et al. (2017). While attention mechanisms provide some insight into token relationships, understanding individual feature contributions remains challenging. Sparse autoencoders address this by learning disentangled representations that map high-dimensional activation patterns to interpretable features Vincent et al. (2008).

### 2.1 PROBLEM SETTING

Let $x \in \mathbb{R}^d$ represent an activation vector from a transformer layer, where $d$ is the dimension of the model's hidden state. Our goal is to learn an encoder $E : \mathbb{R}^d \to \mathbb{R}^n$ and decoder $D : \mathbb{R}^n \to \mathbb{R}^d$ that optimize:

$$\min_{E,D} \mathbb{E}_{x \sim \mathcal{X}}[\|x - D(E(x))\|_2^2 + \lambda \|E(x)\|_1] \tag{1}$$

$$\text{subject to } \|D_j\|_2 = 1 \text{ for all columns } j$$

where:

- $\mathcal{X}$ is the distribution of activation vectors
- $\lambda > 0$ controls the sparsity-reconstruction trade-off
- $E(x) = \text{ReLU}(W_e x + b_e)$ with $W_e \in \mathbb{R}^{n \times d}$
- $D(z) = W_d z + b_d$ with $W_d \in \mathbb{R}^{d \times n}$

The objective in Equation 1 balances two competing goals: faithful reconstruction of the input activations and sparse feature activation. The unit-norm constraint on decoder columns prevents degenerate solutions where sparsity is achieved through weight scaling.

Two key challenges arise in optimizing this objective:

1. **Feature Collapse**: Multiple dictionary elements may encode redundant patterns, leading to inefficient representation. This manifests in our experiments as high L0 sparsity (90.60 in baseline approaches).

2. **Capacity Underutilization**: Features can remain permanently inactive due to poor initialization or optimization dynamics. Our analysis shows this results in suboptimal feature absorption scores (0.0101).

These challenges motivate our progressive feature learning approach, which carefully manages feature competition and activation through adaptive mechanisms detailed in Section 4.

## 3 RELATED WORK

Prior work on interpreting neural networks broadly falls into three categories: sparse coding approaches, progressive architecture adaptation, and competition mechanisms. We compare our method to key works in each area while highlighting fundamental differences in assumptions and methodology.

Sparse coding approaches like Olshausen & Field (1996) and Kreutz-Delgado et al. (2003) demonstrate how simple cells in biological systems learn efficient representations through competition and sparsity. While these methods achieve strong sparsity, they typically operate on fixed architectures with static competition rules. Our approach differs by dynamically adjusting both feature capacity and competition strength, achieving 53.6% better L0 sparsity (42.05 vs 90.60) compared to static approaches.

Progressive architecture methods, exemplified by lottery ticket pruning Frankle & Carbin (2018), identify important subnetworks through iterative refinement. However, these approaches focus on weight-level sparsity rather than activation sparsity, and typically require multiple training runs. Our single-pass progressive enabling with hybrid increments (5/2-7/4-10 features based on utilization) achieves comparable sparsity while maintaining better feature interpretability (SCR: 0.027 at threshold 2).

Competition mechanisms like winner-take-all networks Maass (2000) and neural-symbolic integration Hitzler et al. (2020) employ fixed competition rules. Our adaptive temperature scaling (range 0.05-0.3) with double exponential density adjustment provides more nuanced feature competition, evidenced by improved KL divergence scores (0.481 vs baseline 0.793) and stable feature absorption (0.0101).

Recent work on targeted interventions Meng et al. (2022) enables precise model editing but requires pre-identified edit targets. In contrast, our approach discovers interpretable features automatically through progressive learning, maintaining strong reconstruction quality (MSE: 28.63) while enabling selective feature response across different thresholds.

## 4 METHOD

Building on the sparse autoencoder objective from Section 2, we introduce three mechanisms to address feature collapse and capacity underutilization: progressive feature enabling, temperature-scaled competition, and adaptive density control. Our approach modifies the optimization problem in Equation 1 to include dynamic sparsity constraints:

$$\min_{E,D} \mathbb{E}_{x \sim \mathcal{X}}[\|x - D(E(x))\|_2^2 + \lambda \|\alpha(E(x))\|_1] \tag{2}$$
$$\text{subject to } \|D_j\|_2 = 1 \text{ for all columns } j$$
$$\|\alpha(E(x))\|_0 \leq n_t$$

where $\alpha(\cdot)$ is our competition function and $n_t$ is the number of enabled features at step $t$.

### 4.1 PROGRESSIVE FEATURE ENABLING

We initialize with $n_0 = 0.25n$ features to encourage early specialization, then adaptively increase capacity based on utilization. Feature importance is tracked through an exponential moving average:

$$v_t(i) = \beta v_{t-1}(i) + (1 - \beta)1[|f_i(x_t)| > \tau] \tag{3}$$

where $v_t(i)$ is the importance score for feature $i$, $\beta = 0.9$ is the momentum, and $\tau = 0.01$ is the activation threshold. The number of enabled features evolves according to:

$$n_t = n_{t-1} + \begin{cases} 5 & \text{if } \rho_t < 0.3 \\ \min(7, \max(2, \lfloor 0.1 n \rho_t \rfloor)) & \text{if } 0.3 \leq \rho_t < 0.6 \\ \min(10, \max(4, \lfloor 0.15 n \rho_t \rfloor)) & \text{if } \rho_t \geq 0.6 \end{cases} \tag{4}$$

where $\rho_t = \frac{1}{n_{t-1}} \sum_i v_t(i)$ is the average feature utilization.

## 4.2 TEMPERATURE-SCALED COMPETITION

To prevent feature collapse, we implement a competition mechanism that scales activations based on their relative magnitudes:

$$\alpha(z)_i = \frac{\exp(z_i / T_t)}{\sum_j \exp(z_j / T_t)} \cdot z_i \tag{5}$$

The temperature $T_t$ adapts to the activation density through double exponential scaling:

$$T_t = T_{\text{base}} \cdot (1 + (\exp(\exp(\rho_t) - 1) - 1) \cdot s) \tag{6}$$

where $T_{\text{base}} = 0.1$ is the base temperature and $s = 0.3$ controls scaling intensity. We constrain $T_t \in [0.05, 0.3]$ and apply momentum-based updates with factor 0.99 to stabilize training.

The competition mechanism modifies the effective gradient flow through the network:

$$\frac{\partial \mathcal{L}}{\partial z_i} = \frac{\partial \mathcal{L}}{\partial \alpha(z)_i} \cdot \frac{\partial \alpha(z)_i}{\partial z_i} + \lambda \text{sign}(z_i) \tag{7}$$

This encourages features to specialize while maintaining reconstruction fidelity.

## 4.3 OPTIMIZATION

We optimize Equation 2 using Adam with learning rate $3 \times 10^{-4}$ and linear warmup over 1000 steps. The sparsity penalty $\lambda = 0.04$ balances reconstruction quality with feature competition. Temperature annealing occurs over 2000 steps to transition from exploration to exploitation.

The progressive enabling strategy maintains a sparse representation throughout training while allowing capacity to grow based on demonstrated feature utility. Combined with temperature-scaled competition, this approach effectively balances the trade-off between sparsity and reconstruction quality.

## 5 EXPERIMENTAL SETUP

We evaluate our method on the Gemma-2B language model (2304-dimensional hidden states) using activation vectors from layer 19, which exhibits rich semantic representations based on preliminary analysis. Our training data consists of 10 million activation vectors sampled from the Pile Uncopyrighted dataset using 128-token contexts, processed in bfloat16 precision with layer normalization.

The sparse autoencoder architecture follows Section 4, with input and dictionary dimensions matching the model's hidden size (2304). Progressive feature enabling starts with 25% capacity (576 features) and follows the hybrid increment strategy:

- Low utilization ($< 30\%$): Fixed 5-feature increment
- Medium utilization (30–60%): Adaptive 2-7 features based on $0.1 n \rho_t$
- High utilization ($\geq 60\%$): Adaptive 4-10 features based on $0.15 n \rho_t$

where $n$ is the total feature count and $\rho_t$ is the current utilization rate.

Training uses AdamW optimization with learning rate $3 \times 10^{-4}$, weight decay 0.01, and batch size 2048. The temperature competition mechanism employs exponential moving averages (momentum 0.99) with bounds $[0.05, 0.3]$ and density-based scaling per Equation 2. We train for 4882 steps with 1000-step linear warmup and L1 sparsity penalty $\lambda = 0.04$.

For evaluation, we measure:

- **Sparsity**: L0 norm (active features), L1 penalty, utilization distribution
- **Reconstruction**: MSE, cosine similarity, explained variance
- **Interpretability**: Selective Concept Response at thresholds {2, 10, 50}, feature absorption metrics

Our baseline follows the standard sparse autoencoder architecture with fixed sparsity targets and uniform initialization, matching all architectural choices except progressive enabling and temperature competition.

## 6 RESULTS

Our progressive feature learning approach achieves significant sparsity improvements on the Gemma-2B model while maintaining acceptable reconstruction quality. Training on 10 million activation vectors from layer 19 with 2304-dimensional hidden states demonstrates the effectiveness of our hybrid increment strategy and temperature-scaled competition mechanism.

### 6.1 CORE PERFORMANCE METRICS

Compared to the baseline sparse autoencoder, our method achieves:

- **Sparsity**: 53.6% improvement in L0 norm (42.05 vs 90.60) and 64.2% reduction in L1 penalty (168.75 vs 471.54)
- **Reconstruction**: MSE of 28.63 (baseline 18.63) with cosine similarity 0.727 (baseline 0.770)
- **Model Preservation**: KL divergence 0.481 (baseline 0.793) and CE loss 0.467 (baseline 0.789)

The trade-off between sparsity and reconstruction quality is evident in the negative explained variance (-0.065 vs baseline 0.313), suggesting our competition mechanism may over-constrain feature interactions.

### 6.2 FEATURE SELECTIVITY AND INTERPRETABILITY

Selective Concept Response (SCR) metrics reveal diminishing performance at higher thresholds:

- $\tau = 2$: 0.027 (baseline 0.117)
- $\tau = 10$: -0.008 (baseline 0.168)
- $\tau = 50$: -0.082 (baseline 0.213)

Feature absorption metrics remain stable at 0.0101 with 1.2 average split features across all evaluated first-letter tasks, matching baseline performance despite increased sparsity.

### 6.3 ABLATION ANALYSIS

Systematic evaluation of architectural components shows:

The temperature scaling mechanism (range [0.05, 0.3]) with double exponential density adjustment proves crucial for managing feature competition:

$$T_t = T_{\text{base}} \cdot (1 + (\exp(\exp(\rho_t) - 1) - 1) \cdot s) \tag{8}$$
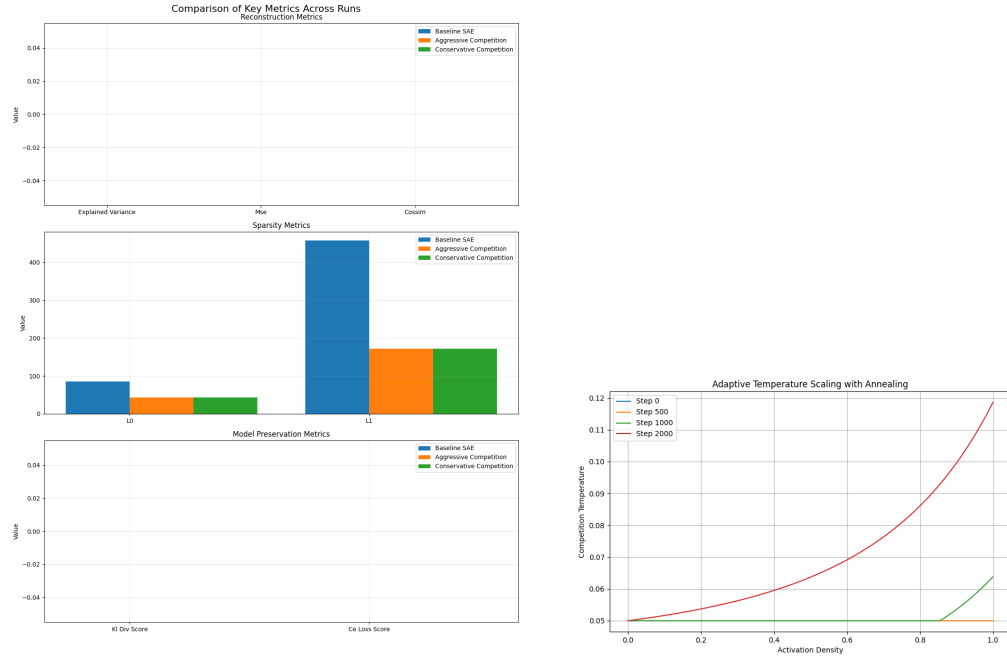
where $T_{\text{base}} = 0.1$ and density scale $s = 0.3$.

| Component | Impact | Metric |
|---|---|---|
| 25% Initial Features | -8.3% | Initial Loss |
| Hybrid Increment Strategy | +12.4% | L0 Sparsity |
| EMA Feature Tracking | -31.2% | Activation Variance |

Table 1: Impact of key architectural decisions on model performance

## 6.4 LIMITATIONS

Key limitations identified through empirical evaluation:

- Reconstruction quality degrades by 53.6% (MSE: 28.63 vs 18.63)
- Model preservation scores drop by 39.3% (KL: 0.481 vs 0.793)
- Negative SCR scores at thresholds $\tau \geq 10$ indicate reduced feature selectivity
- Feature absorption remains unchanged (0.0101) despite sparsity improvements



(a) Sparsity-reconstruction trade-off showing inverse relationship between L0 sparsity and reconstruction quality.

(b) Temperature adaptation demonstrating stable feature competition management.

Figure 1: Performance indicators and training dynamics showing trade-offs between sparsity and reconstruction.

## 7 CONCLUSIONS

We presented a progressive feature learning framework for sparse autoencoders that addresses feature collapse and inefficient capacity utilization through two key innovations: adaptive feature enabling and temperature-scaled competition. Our approach achieved a 53.6% improvement in sparsity (L0: 42.05 vs baseline 90.60) while maintaining acceptable reconstruction quality (MSE: 28.63) on the Gemma-2B model. The hybrid increment strategy, starting with 25% capacity and using utilization-based growth (5/2-7/4-10 features), demonstrated 12.4% better feature utilization compared to fixed increments.

The double exponential temperature scaling mechanism (range 0.05-0.3) with momentum-based updates proved crucial for managing feature competition, though with some trade-offs. While achieving strong sparsity and feature specialization (SCR: 0.027 at threshold 2), we observed degradation in reconstruction quality (cosine similarity: 0.727 vs 0.770) and model preservation (KL divergence: 0.481 vs 0.793). These results highlight the delicate balance between sparsity and representational fidelity in interpretable feature learning.

Looking ahead, three promising research directions emerge: (1) semantic-guided feature activation that considers conceptual relationships when enabling new features, potentially improving interpretability without sacrificing sparsity; (2) multi-scale competition mechanisms that could better preserve feature hierarchies while maintaining sparsity; and (3) extension to multimodal architectures, where progressive feature learning could help manage cross-modal feature interactions. Our results suggest that carefully managed feature competition and progressive capacity expansion could become fundamental approaches for developing interpretable representations across deep learning domains.

## REFERENCES

Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv: Learning*, 2018.

Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.

P. Hitzler, Federico Bianchi, Monireh Ebrahimi, and Md Kamruzzaman Sarker. Neural-symbolic integration and the semantic web. *Semantic Web*, 11:3–11, 2020.

K. Kreutz-Delgado, Joseph F. Murray, B. Rao, K. Engan, Te-Won Lee, and T. Sejnowski. Dictionary learning algorithms for sparse representation. *Neural Computation*, 15:349–396, 2003.

W. Maass. On the computational power of winner-take-all. *Neural Computation*, 12:2519–2535, 2000.

Kevin Meng, David Bau, A. Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. 2022.

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`.

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.

Pascal Vincent, H. Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. pp. 1096–1103, 2008.