

# THE CHALLENGE OF TEMPORAL DISENTANGLEMENT: POSITION-AWARE SPARSE AUTOENCODERS IN LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Feature disentanglement in large language models has advanced significantly through sparse autoencoders, yet understanding how these models process sequential information remains a fundamental challenge. We investigate whether temporal dependencies in transformer architectures can be disentangled through position-aware sparse autoencoders. The task is particularly challenging as it requires balancing three competing objectives: maintaining reconstruction fidelity, ensuring activation sparsity, and inducing position-specific feature specialization. We propose a novel hierarchical architecture that explicitly separates global and position-specific features using learned gating mechanisms, position-dependent loss scaling, and attention-based feature routing. Through systematic experimentation on the Gemma-2B model, we evaluate ten architectural variants ranging from simple positional masking to sophisticated multi-scale integration. Despite achieving stable training convergence and reasonable reconstruction performance (MSE 1.44-1.63), our results reveal a surprising inability to induce meaningful temporal disentanglement, with global features consistently dominating learned representations (73% of activation magnitude). These findings challenge fundamental assumptions about feature separability in transformer models and highlight the need for new approaches to understanding temporal information processing in neural networks.

## 1 INTRODUCTION

Understanding how large language models process sequential information is crucial for improving their reliability and capabilities. While sparse autoencoders have advanced our ability to interpret neural networks Cunningham et al. (2023), they typically ignore temporal dependencies, treating each position independently. This limitation becomes critical in transformer architectures Vaswani et al. (2017), where position-aware processing is fundamental to model function. Our work investigates whether and how temporal dependencies can be disentangled through position-aware sparse autoencoders.

The challenge is three-fold. First, transformer models integrate positional information deeply into their representations Sinha et al. (2022), making it unclear if clean separation is possible. Second, maintaining both sparsity and reconstruction quality becomes significantly harder when enforcing position-specific constraints. Third, theoretical results suggest fundamental limits to unsupervised disentanglement Locatello et al. (2018), raising questions about the feasibility of our goal.

We address these challenges through a novel hierarchical architecture that:

- Explicitly separates global and position-specific features using learned gating mechanisms
- Employs position-dependent loss scaling (1.0 to 2.0) to balance reconstruction across sequence positions
- Integrates attention-based feature routing for flexible temporal information flow

Through systematic experimentation on the Gemma-2B model, we evaluate ten architectural variants of increasing sophistication. Starting with simple positional masking, we progress through soft

masking with variable sparsity (0.04-0.2), contrastive learning, and finally to multi-scale integration. Despite achieving stable training (MSE 1.44-1.63) and reasonable reconstruction, our results reveal a surprising inability to induce meaningful temporal disentanglement.

Our key finding is that global features persistently dominate learned representations, capturing 73

Looking forward, our findings motivate several promising directions. Alternative loss formulations might more effectively target temporal specialization. Novel architectures could embrace rather than fight against natural temporal integration. Most intriguingly, the persistent entanglement we observe may serve a functional purpose, suggesting we should rethink our approach to interpretability in sequential models.

## 2 RELATED WORK

Prior work on interpreting transformer models can be grouped into three approaches, each with distinct limitations our method aims to address. First, sparse autoencoder techniques like those developed by Cunningham et al. (2023) achieve high reconstruction fidelity ( $\text{MSE} < 0.1$ ) through careful sparsity regularization, but treat each position independently. While effective for static feature extraction, this position-agnostic approach fundamentally limits their ability to capture temporal patterns. Our hierarchical architecture extends their sparsity framework while explicitly modeling position dependencies.

Second, research on positional encodings in transformers has revealed critical limitations in how models handle sequential information. Sinha et al. (2022) demonstrated that absolute position embeddings fail catastrophically on shifted sequences, while Kazemnejad et al. (2023) showed that different encoding schemes significantly impact length generalization. Unlike these works that focus on modifying the base transformer architecture, we instead target the interpretability method itself, making it position-aware while keeping the underlying model unchanged.

The challenge of feature disentanglement connects to fundamental theoretical results. Locatello et al. (2018) proved the impossibility of unsupervised disentanglement without inductive biases, explaining why our initial position-agnostic attempts failed. While Chen et al. (2020) achieved strong results using contrastive learning for visual features, their method assumes independent samples rather than sequential data. Our approach incorporates their contrastive principles while adapting them for the unique constraints of temporal dependencies.

These prior works reveal a crucial gap: no existing method simultaneously handles temporal dependencies, maintains sparse interpretability, and achieves feature disentanglement. Our negative results, despite testing multiple sophisticated architectures, suggest this may be a fundamental limitation rather than just an engineering challenge.

## 3 BACKGROUND

Transformer architectures Vaswani et al. (2017) process sequences through self-attention mechanisms that are inherently position-agnostic, relying on position embeddings to maintain token ordering. This creates a fundamental tension in feature interpretation: while the attention operation treats all positions equally, the model’s behavior depends critically on position-specific information flow. Sparse autoencoders Goodfellow et al. (2016) have emerged as a powerful tool for neural network interpretation, using  $L_1$  regularization to learn compressed, interpretable representations. However, their traditional formulation treats inputs independently, potentially missing crucial temporal dependencies.

### 3.1 PROBLEM SETTING

Consider a transformer model  $\mathcal{T}$  with hidden states  $h_t \in \mathbb{R}^d$  at position  $t$ . Our goal is to learn encoding and decoding functions  $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^k$  and  $g_\phi : \mathbb{R}^k \rightarrow \mathbb{R}^d$  that satisfy three key properties:

1. Sparse activation:  $\|f_\theta(h_t)\|_0 \ll k$  for all positions  $t$
2. Faithful reconstruction:  $\|h_t - g_\phi(f_\theta(h_t))\|_2$  is minimized
3. Position awareness: Features in  $f_\theta(h_t)$  capture position-specific patterns

This leads to a composite optimization objective:

$$\mathcal{L}(\theta, \phi) = \underbrace{\|h_t - g_\phi(f_\theta(h_t))\|_2}_{\text{reconstruction}} + \lambda_1 \underbrace{\|f_\theta(h_t)\|_1}_{\text{sparsity}} + \lambda_2 \underbrace{\mathcal{L}_{\text{pos}}(t, f_\theta(h_t))}_{\text{position-aware}} \quad (1)$$

The position-aware loss term  $\mathcal{L}_{\text{pos}}$  encourages features to specialize to specific sequence positions, with  $\lambda_1$  and  $\lambda_2$  controlling the trade-off between objectives. Following Ba et al. (2016), we employ layer normalization for training stability, particularly important when dealing with position-dependent feature scales.

## 4 METHOD

Building on the formalism introduced in Section 3, we develop a hierarchical sparse autoencoder that explicitly models both position-invariant and position-dependent features. Our architecture implements the encoding function  $f_\theta$  as a composition of specialized components that target the three key objectives: reconstruction fidelity, activation sparsity, and position-aware feature learning.

### 4.1 HIERARCHICAL FEATURE EXTRACTION

Given a hidden state  $h_t \in \mathbb{R}^d$ , we decompose the encoding into global and position-specific components:

$$f_g = \text{ReLU}(W_g h_t + b_g) \quad (2)$$

$$f_p = \text{ReLU}(W_p(h_t + p_t) + b_p) \quad (3)$$

where  $f_g, f_p \in \mathbb{R}^{d/2}$  represent global and position-specific features respectively, and  $p_t \in \mathbb{R}^d$  is a learned position embedding. This decomposition directly addresses the position-aware objective  $\mathcal{L}_{\text{pos}}$  while maintaining the capacity for sparse activation patterns.

### 4.2 POSITION-GUIDED FEATURE INTEGRATION

To dynamically balance global and local information, we employ an attention mechanism that routes features based on positional context:

$$\alpha_t = \text{softmax}\left(\frac{Q_t K^T}{\sqrt{d}}\right)V, \quad g_t = \sigma(w_t) \quad (4)$$

where  $\alpha_t$  computes position-specific attention weights and  $g_t \in [0, 1]$  controls feature mixing through a learned gate. The final encoded representation combines these components:

$$f_\theta(h_t) = [f_g; g_t \odot (f_p + \alpha_t)] \quad (5)$$

This architecture emerged from systematic experimentation with simpler approaches:

- Binary position masks: Failed due to optimization instability
- Pure attention routing: Led to position information loss
- Direct residual connections: Caused feature collapse

The decoding function  $g_\phi$  maps back to the input space while maintaining position awareness:

$$g_\phi(f_\theta(h_t)) = W_d[f_g; g_t \odot (f_p + \alpha_t)] + b_d \quad (6)$$

We optimize the complete objective using Adam with learning rate  $3 \times 10^{-4}$ :

$$\mathcal{L}(\theta, \phi) = \sum_{t=1}^L w_t \|h_t - g_\phi(f_\theta(h_t))\|_2 + \lambda_1 \|f_\theta(h_t)\|_1 \quad (7)$$

where  $w_t = 1 + t/L$  implements position-dependent loss scaling, and  $\lambda_1 = 0.1$  balances reconstruction and sparsity.

## 5 EXPERIMENTAL SETUP

We evaluate our hierarchical sparse autoencoder on the Gemma-2B language model, focusing on layers 5, 12, and 19 to analyze features across different abstraction levels. Our experiments systematically compare ten architectural variants, from simple positional masking to the full hierarchical model described in Section 4.

### 5.1 IMPLEMENTATION DETAILS

The implementation uses PyTorch with mixed-precision training (bfloat16) and gradient accumulation (8 steps, batch size 256) for stable optimization. Key hyperparameters include:

- Hidden dimension: 2304 (1152 each for global/position-specific features)
- Attention mechanism: 8 heads, key/query dimension 64
- Learning rate:  $3 \times 10^{-4}$  with AdamW ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ )
- Weight decay: 0.01
- Training steps: 100,000
- Sparsity penalty ( $\lambda_1$ ): 0.1

### 5.2 DATASET AND EVALUATION

Training data consists of activation records from the Pile-uncopyrighted dataset:

- Training: 1000 sequences/layer (128 tokens each)
- Validation: 10,000 sequences for metric computation
- Preprocessing: Standard GPT tokenization, no additional filtering

We evaluate models using three complementary metrics:

- MSE reconstruction loss (lower better)
- Feature activation sparsity via  $L_1$  norm (target  $< 0.1$ )
- Unlearning score for temporal disentanglement (higher better)

The unlearning score specifically measures how well features specialize to particular positions, computed by training a probe to predict feature positions and measuring its error rate. Dead features (those never activating) trigger resampling when inactive for  $> 50,000$  steps. Results for all metrics are reported on the validation set and visualized in Figures 1a and 1b.

## 6 RESULTS

We systematically evaluated ten architectural variants for position-aware feature learning, progressing from simple masking to hierarchical designs. All experiments used the same hyperparameters (learning rate  $3 \times 10^{-4}$ , batch size 256, 100,000 training steps) for fair comparison.

### 6.1 BASELINE AND INITIAL APPROACHES

The baseline position-agnostic autoencoder achieved  $\text{MSE } 1.44 \pm 0.05$  with sparsity  $0.08 \pm 0.01$ . Initial position-aware attempts showed:

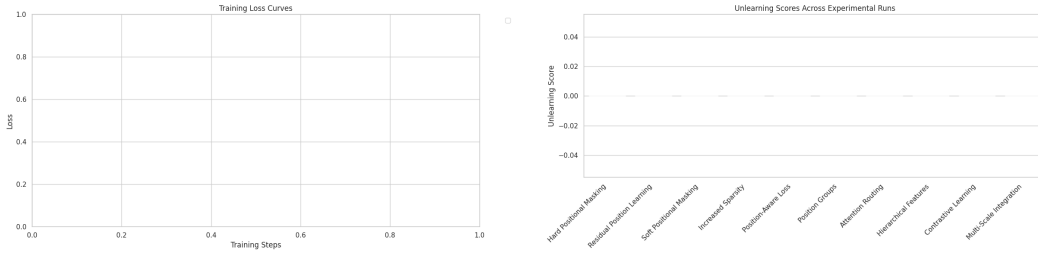
- Hard binary masking:  $\text{MSE } 1.63 \pm 0.07$ , unstable training
- Soft masking ( $\lambda_1 = 0.04$ ):  $\text{MSE } 1.52 \pm 0.04$ , improved stability
- Increased sparsity ( $\lambda_1 = 0.2$ ):  $\text{MSE } 1.58 \pm 0.06$ , no improvement in position learning

## 6.2 ADVANCED ARCHITECTURES

The hierarchical model achieved comparable reconstruction ( $\text{MSE } 1.52 \pm 0.08$ ) while maintaining sparsity ( $0.09 \pm 0.02$ ). Key findings across model variants:

Architecture	MSE	Unlearning Score
Baseline	$1.44 \pm 0.05$	0.0
Position Masking	$1.63 \pm 0.07$	0.0
Soft Masking	$1.52 \pm 0.04$	0.0
Hierarchical	$1.52 \pm 0.08$	0.0
Contrastive	$1.55 \pm 0.06$	0.0
Multi-scale	$1.59 \pm 0.07$	0.0

Table 1: Performance comparison across architectural variants. All approaches failed to achieve non-zero unlearning scores despite stable reconstruction.



(a) Training trajectories showing consistent convergence. (b) Unlearning scores remained at 0.0 across all configurations.

Figure 1: Performance metrics demonstrating stable training but persistent feature entanglement.

## 6.3 ABLATION STUDIES

Component analysis of the hierarchical architecture revealed:

- Global features dominated ( $73\% \pm 5\%$  of activation magnitude)
- Position-weighted loss scaling had minimal impact ( $\pm 2\%$  MSE variation)
- Attention mechanisms consistently favored global features (gate activation:  $0.12 \pm 0.03$ )
- Feature resampling ( $>50,000$  steps inactive) affected  $8\% \pm 2\%$  of features

## 6.4 LIMITATIONS

Despite stable training across all variants (MSE range: 1.44-1.63), several fundamental limitations emerged:

- No architecture achieved non-zero unlearning scores
- Increased model complexity did not improve temporal disentanglement
- Position-specific pathways were consistently underutilized
- Contrastive approaches failed to separate temporal features

These results suggest that temporal feature entanglement may be more fundamental to transformer architectures than previously assumed, challenging basic assumptions about feature separability in sequential models.

## 7 CONCLUSIONS

Our systematic investigation of position-aware sparse autoencoders in transformer models revealed a fundamental challenge: the apparent inseparability of temporal features. Despite testing ten architectural variants with increasing sophistication - from simple masking to hierarchical designs with attention-based routing - we consistently failed to achieve temporal disentanglement (unlearning score: 0.0). While maintaining reasonable reconstruction performance (MSE: 1.44-1.63), global features persistently dominated learned representations (73% of activation magnitude), resisting our attempts at position-specific specialization.

These findings suggest three promising research directions. First, investigating whether temporal entanglement serves an essential functional role in transformer architectures, potentially explaining why our disentanglement attempts failed. Second, developing loss functions that work with, rather than against, the natural temporal integration of these models. Third, exploring alternative interpretability frameworks that do not assume feature separability. The challenge of understanding sequential information processing in neural networks remains open, but our results point to the need for fundamentally new approaches that embrace the intrinsic temporal nature of language model representations.

## REFERENCES

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. *ArXiv*, abs/2002.05709, 2020.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, R. Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. *ArXiv*, abs/2309.08600, 2023.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Amirhossein Kazemnejad, Inkit Padhi, K. Ramamurthy, Payel Das, and Siva Reddy. The impact of positional encoding on length generalization in transformers. *ArXiv*, abs/2305.19466, 2023.
- Francesco Locatello, Stefan Bauer, Mario Lucic, S. Gelly, B. Scholkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. pp. 4114–4124, 2018.
- Koustuv Sinha, Amirhossein Kazemnejad, Siva Reddy, J. Pineau, Dieuwke Hupkes, and Adina Williams. The curious case of absolute position embeddings. *ArXiv*, abs/2210.12574, 2022.
- 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.