

AbsTopK: Rethinking Sparse Autoencoders For Bidirectional Features

Xudong Zhu Mohammad Mahdi Khalili Zihui Zhu

Department of Computer Science and Engineering
The Ohio State University

ICLR 2026
April 23–27, 2026

Mechanistic Interpretability

- 1 Mechanistic interpretability aims to decompose hidden representations into meaningful features.
- 2 Sparse autoencoders (SAEs) aim to recover these features through a sparse latent representation.

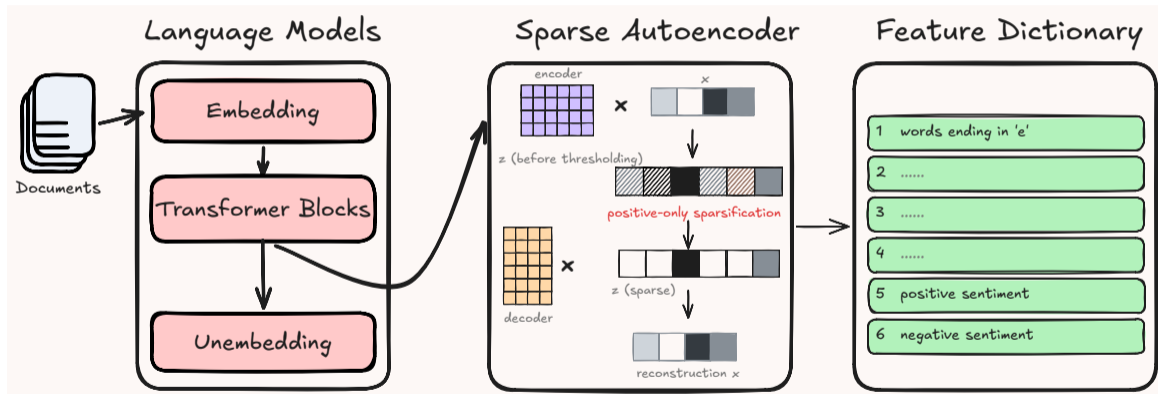


Figure: Overall pipeline of a sparse autoencoder.

Sparse Autoencoders (SAEs)

Assume the hidden representation \mathbf{x} can be decomposed as

$$\mathbf{x} = \sum_{p=1}^P \alpha_p \mathbf{h}_p + \text{residual}.$$

An SAE maps the input \mathbf{x} to a sparse code \mathbf{z} , and then reconstructs the original representation:

$$\text{encoder: } \mathbf{z} = \sigma(\mathbf{W}^T \mathbf{x} + \mathbf{b}_e), \quad \text{decoder: } \hat{\mathbf{x}} = \mathbf{D}\mathbf{z} + \mathbf{b}.$$

It is trained by minimizing the following objective:

$$\min_{\substack{\mathbf{D}, \mathbf{W} \in \mathbb{R}^{d \times P} \\ \mathbf{b} \in \mathbb{R}^d, \mathbf{b}_e \in \mathbb{R}^P}} \mathbb{E}_{\mathbf{x}} \left[\frac{1}{2} \|\mathbf{x} - (\mathbf{D}\mathbf{z} + \mathbf{b})\|_2^2 + \lambda R(\mathbf{z}) \right].$$

Concepts

Each row of the decoder \mathbf{D} can be interpreted as a concept in the learned dictionary, so the reconstructed hidden state is a linear combination of these concepts:

$$\hat{\mathbf{x}} = \sum_{p=1}^P z_p \mathbf{d}_p + \mathbf{b}.$$

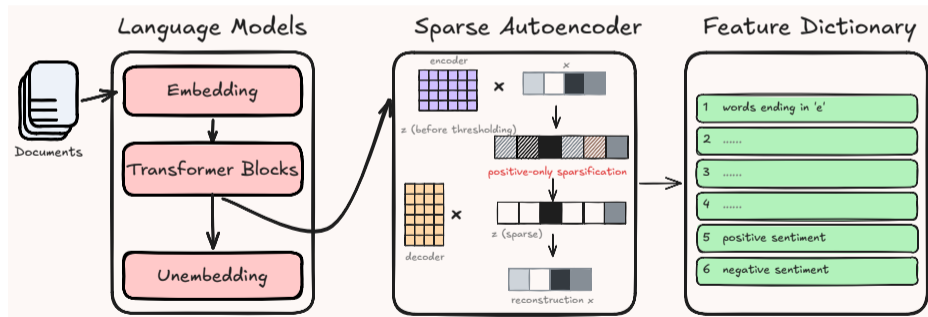


Figure: Each decoder vector can be viewed as a concept in the learned feature dictionary.

The Proximal Perspective on SAEs

What distinguishes different SAE variants?

Our key observation is that several popular SAEs admit a unified view through proximal optimization, rooted in classical sparse coding and dictionary learning.

$$\min_{\substack{\mathbf{D}, \mathbf{W} \in \mathbb{R}^{d \times P} \\ \mathbf{b} \in \mathbb{R}^d, \mathbf{b}_e \in \mathbb{R}^P}} \mathbb{E}_{\mathbf{x}} \left[\frac{1}{2} \|\mathbf{x} - (\mathbf{D}\mathbf{z} + \mathbf{b})\|_2^2 + \lambda R(\mathbf{z}) \right], \quad \text{where } \mathbf{z} = \text{prox}_{\lambda R}(\mathbf{W}^\top \mathbf{x} + \mathbf{b}_e).$$

- ReLU SAE:

$$R(\mathbf{z}) = \|\mathbf{z}\|_1 + \iota_{\{\mathbf{z} \geq 0\}}(\mathbf{z}) \implies \text{prox}_{\lambda R} = \text{ReLU}_\lambda$$

- JumpReLU SAE:

$$R(\mathbf{z}) = \|\mathbf{z}\|_0 + \iota_{\{\mathbf{z} \geq 0\}}(\mathbf{z}) \implies \text{prox}_{\lambda R} = \text{JumpReLU}_{\sqrt{2\lambda}}$$

- TopK SAE:

$$R(\mathbf{z}) = \iota_{\{\|\mathbf{z}\|_0 \leq k, \mathbf{z} \geq 0\}}(\mathbf{z}) \implies \text{prox}_{\lambda R} = \text{TopK}_k$$

The Non-Negativity Constraint

Popular SAE variants all enforce a non-negativity constraint:

$$\text{ReLU SAE: } R(\mathbf{z}) = \|\mathbf{z}\|_1 + \iota_{\{\mathbf{z} \geq 0\}}(\mathbf{z})$$

$$\text{JumpReLU SAE: } R(\mathbf{z}) = \|\mathbf{z}\|_0 + \iota_{\{\mathbf{z} \geq 0\}}(\mathbf{z})$$

$$\text{TopK SAE: } R(\mathbf{z}) = \iota_{\{\|\mathbf{z}\|_0 \leq k, \mathbf{z} \geq 0\}}(\mathbf{z})$$

- A latent can only activate in one direction.
- Opposite concepts are split into different features.
- AbsTopK removes this constraint and allows bidirectional features.

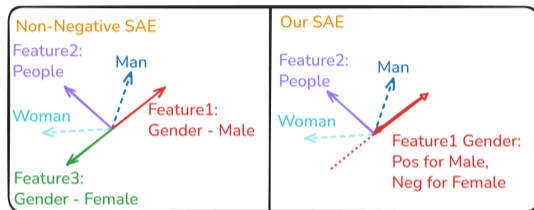


Figure: Non-negativity fragments a semantic axis into separate features, while AbsTopK captures both directions in one latent.

Beyond Non-Negativity: AbsTopK

Specifically, we use the regularizer $R(\mathbf{z}) = \iota_{\{\|\mathbf{z}\|_0 \leq k\}}$ which removes the non-negativity constraint present in the TopK-inducing regularizer. The corresponding proximal operator is

$$\text{prox}_{\lambda R}(\mathbf{u}) = \arg \min_{\mathbf{z} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{u} - \mathbf{z}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{z}\|_0 \leq k, \quad (1)$$

whose closed-form solution is further given by

$$(\text{prox}_R(\mathbf{u}))_i = (\text{AbsTopK}_k(\mathbf{u}))_i = \begin{cases} u_i, & i \in \mathcal{H}_k(\mathbf{u}), \\ 0, & i \notin \mathcal{H}_k(\mathbf{u}), \end{cases} \quad (2)$$

where \mathcal{H}_k denotes the indices of the k largest (in modulus) components.

Beyond Non-Negativity: AbsTopK

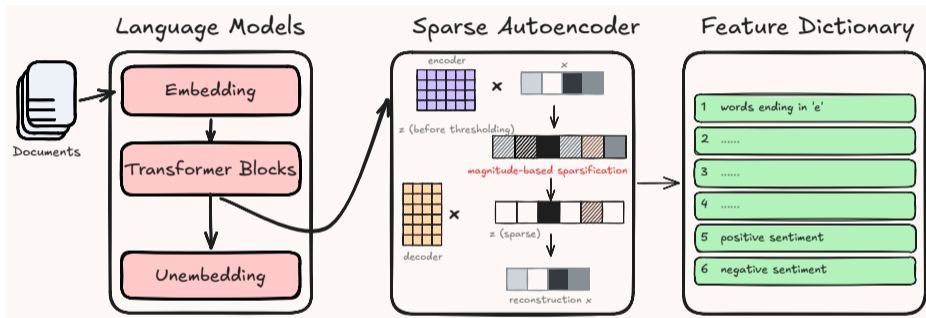


Figure: Instead of keeping only positive activations, AbsTopK selects coefficients by magnitude, preserving both directions of a semantic axis.

AbsTopK Learns Bidirectional Features

SAE feature 15741: Positive activation on **man**

The **man** who had traveled across ...

Despite the storm outside , the **man** remained calm and focused on her work ...

In the bustling market , a **man** stood by the flower stall ...

SAE feature 15741: Negative activation on **woman**

The **woman** who had traveled across ...

Despite the storm outside , the **woman** remained calm and focused on her work ...

In the bustling market , a **woman** stood by the flower stall ...

Figure: A single AbsTopK feature captures both directions of a semantic axis: positive activation aligns with man, while negative activation aligns with woman.

AbsTopK Learns Better SAEs

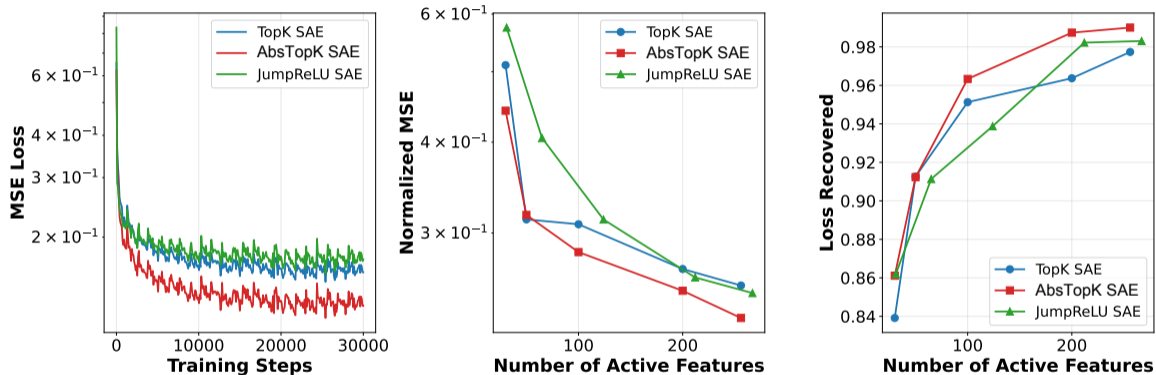


Figure: AbsTopK achieves lower reconstruction error and higher loss recovered than TopK and JumpReLU across different sparsity levels.

AbsTopK Improves Steering Performance

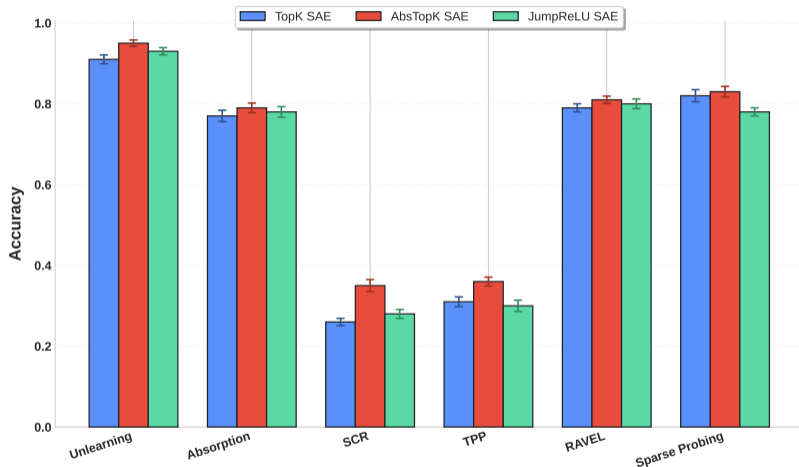


Figure: AbsTopK achieves the best or near-best accuracy across multiple downstream tasks.

- 1 Popular SAEs can be understood through a unified proximal perspective, that their shared non-negativity constraint fragments semantic axes.
- 2 AbsTopK removes this limitation to learn bidirectional features.
- 3 Empirically, this leads to better interpretability, better reconstruction, and stronger downstream performance.