# Does higher interpretability imply better utility? A Pairwise Analysis on Sparse Autoencoders

## **Anonymous Author(s)**

Affiliation Address email

## **Abstract**

Sparse Autoencoders (SAEs) are widely used to steer large language models (LLMs), based on the assumption that their interpretable features naturally enable effective model behavior steering. Yet, a fundamental question remains unanswered: does higher interpretability indeed imply better steering utility? To answer this question, we train 90 SAEs across three LLMs (Gemma-2-2B, Qwen-2.5-3B, Gemma-2-9B), spanning five architectures and six sparsity levels, and evaluate their interpretability and steering utility based on SAEBENCH [Karvonen et al., 2025] and AXBENCH [Wu et al., 2025] respectively, and perform a rankagreement analysis via Kendall's rank coefficients  $\tau_b$ . Based on the framework, Our analysis reveals only a relatively weak positive association ( $\tau_b \approx 0.298$ ), indicating that interpretability is an insufficient proxy for steering performance. We conjecture the interpretability-utility gap may stem from the selection of SAE features as not all of them are equally effective for steering. To further find features that truly steer the behavior of LLMs, we propose a novel selection criterion:  $\Delta$ Token Confidence, which measures how much amplifying a feature changes the next token distribution. We show that our method improves the steering performance of three LLMs by 52.52% compared to the current best output score-based criterion [Arad et al., 2025]. Strikingly, after selecting features with high  $\Delta$  Token Confidence, the correlation between interpretability and utility vanishes  $(\tau_h \approx 0)$ , and can even become negative. This further highlights the divergence between interpretability and utility for the most effective steering features.

#### 22 1 Introduction

2

3

5

6

7

9

10

11

12

13

14

15

16

17

18

19

20

21

23

24

27

28

29

30

31

34

35

As Large Language Models (LLMs) become more widely used in real-world applications, ensuring the safety of their outputs is increasingly important [Kumar et al., 2024, Ji et al., 2023, Inan et al., 2023]. Reliable and controllable behavior is essential for deploying these LLMs in more situations [Chen et al., 2024]. Fine-tuning is the standard way to improve controllability, but it requires labeled data, significant training time, and compute resources [Hu et al., 2022, Wang et al., 2025a]. This has spawned a series of representation-based interventions, i.e., steering, that guide LLM inference by manipulating internal representations, aiming for faster and more lightweight output control [Turner et al., 2023, Turner et al., 2024, Wang et al., 2025b, Stolfo et al., 2025]. However, activation-level edits are often coarse: they mix multiple semantics, a phenomenon called polysemanticity [Bricken et al., 2023]. Recently, Sparse Autoencoders (SAEs) have become a valuable tool in the interpretability field. They are trained to actively decompose the hidden states of the LLM into sparse and human readable features [Templeton et al., 2024, Mudide et al., 2025]. Their interpretable nature has subsequently spurred research into leveraging SAE features for more precise, concept-level control over model behavior [Ferrando et al., 2025, Chalnev et al., 2024].

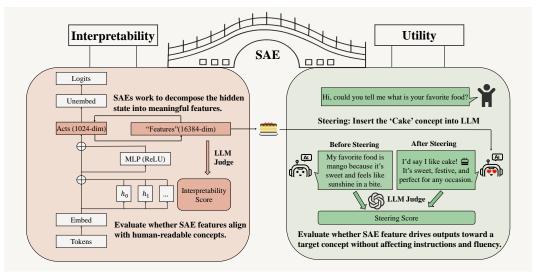


Figure 1: Overview of our goal: building a bridge for SAE interpretability and utility. Interpretability (left): an SAE attached to the LLM decomposes hidden states into sparse, human-describable features. An LLM judge yields an *interpretability score* for the SAE [Paulo et al., 2025]. Utility (right): at inference, we activate a target SAE feature (e.g., 'cake') to steer generation. An LLM judge yields *steering utility score* [Wu et al., 2025].

Despite this progress, a critical question remains unanswered: *does higher interpretability truly imply better utility?* Since SAEs are trained to balance reconstruction and sparsity to yield human-readable features [Cunningham et al., 2023, Makelov, 2024, O'Brien et al.], their utility for down-stream tasks is not a primary objective. Understanding and characterizing this gap is critical to enabling more interpretable and effective steering over the LLM. To this end, we conduct a system-atic study to build a bridge between SAE interpretability and steering utility (see Figure 1).

To perform a comprehensive association analysis, we train 90 SAEs across three LLMs (Gemma-2-2B [Team et al., 2024], Qwen-2.5-3B [Yang et al., 2024], and Gemma-2-9B) spanning diverse architectures and sparsity levels. We compute interpretability using SAEBENCH [Karvonen et al., 2025] and steering utility using AxBENCH [Wu et al., 2025]. Then, we leverage a pairwise-controlled framework to evaluate whether interpretability predicts steering performance across the pool of trained SAEs. To quantify this relationship, we follow the idea of prior works [Jiang et al., 2020, Hu et al., 2024] and measure rank agreement between interpretability and utility using Kendall's rank coefficient  $\tau_b$ . We control confounders with an axis-conditioned analysis, isolating each design axis (architecture, sparsity, model) by varying one at a time and aggregating per-axis metrics.

Furthermore, as identified in Arad et al. [2025], Wu et al. [2025], not all interpretable features in SAE are equally effective for steering. This motivates our next objective to identify the specific features critical for behavior control and steering utility analysis. Motivated by the recent progress on the entropy mechanism in LLM reasoning [Fu et al., 2025, Wang et al., 2025c], we propose an innovative selection criterion for SAE features:  $\Delta$  *Token Confidence*, which measures the degree to which amplifying a single feature shifts the model's next-token distribution. Features that induce the most substantial change in model confidence are identified as high-utility candidates features for steering, as they exert a measurable and targeted influence on model behavior. Finally, we leverage these critical features to conduct a refined analysis of the interpretability-utility gap.

The primary contributions and insights of this paper are summarized as follows:

- 1. (§3.4) Interpretability shows a relatively weak positive association with utility. Across 90 SAEs that are trained across three model sizes, five architectures, and six sparsity levels, we find that a higher *interpretability score* tends to shows a relatively weak positive association with steering performance (the Kendall's rank coefficient  $\tau_b \approx 0.298$ )). This identifies a notable interpretability-utility gap of the existing SAEs.
- 2. (§4.2)  $\triangle$  *Token Confidence* effectively selects features with strong steering performance. To identify the SAE features that are critical for steering, we introduce  $\triangle$  *Token Confidence*, an

innovative metric that identifies steering-critical SAE features by measuring their impact on the model's next-token distribution. When benchmarked against the best existing output score-based method [Arad et al., 2025], our approach yields a substantial 52.52% average improvement in steering score. This result validates the superiority of our method and underscores the critical role of feature selection in characterizing and enhancing the steering utility of SAEs.

3. (§4.3) The interpretability-utility gap widens among high-utility features. By reapplying our association analysis exclusively to SAE features with strong steering utility, we uncover a counterintuitive finding: the interpretability-utility correlation vanishes or even becomes negative (Kendall's rank coefficient  $\tau_b \approx 0$ ). This indicates that for the most effective steering features, interpretability is at best irrelevant and potentially detrimental, further emphasizing the critical nature of the interpretability-utility gap.

Our results demonstrate a significant divergence between SAE's interpretability and steering utility, suggesting that prioritizing interpretability does not enable improved steering performance. This gap highlights a crucial research direction: mitigating it will likely necessitate advanced post-training feature selection protocols or fundamentally new, utility-oriented SAE training paradigms.

## 84 2 Preliminary

74

75

76

77

78

79

#### 85 2.1 Sparse Autoencoders

Sparse Autoencoders (SAEs) decompose internal model activations x into sparse, higherdimensional features h that can be linearly decoded back to the original space [Cunningham et al.,
2023, Leask et al., 2025]. A standard SAE with column-normalized decoder weights [Bricken et al.,
2023, Karvonen et al., 2024] is defined by the following forward map and optimization objective:

$$\mathcal{L} = \|x - \hat{x}\|_{2}^{2} + \lambda \|h\|_{1}$$
, where  $h = \text{ReLU}(W_{E}x + b_{E})$ ,  $\hat{x} = W_{D}h + b_{D}$ ,

where  $W_E, b_E$  are encoder parameters,  $W_D, b_D$  are decoder parameters,  $\hat{x}$  is the reconstruction, and  $\lambda$  controls sparsity. This training balances reconstruction accuracy with sparse representations.

## 92 2.2 Interpretability: Automated Interpretability Score

SAEBENCH [Karvonen et al., 2025] uses an LLM-as-judge [Paulo et al., 2025] to assess each latent: the judge drafts the description from examples and then predicts, on a held-out set, which sequences activate it. The *Automated Interpretability Score* is the average precision of the judge's prediction.

AutoInterp Score = 
$$\frac{1}{M} \sum_{m=1}^{M} \mathbf{1}[\hat{y}_m = y_m],$$

where  $y_m \in \{0,1\}$  indicates whether the latent activates in the sequence m and  $\hat{y}_m$  is the judge's prediction. We use this score as our interpretability metric. For the complete details, see Appendix B.

## 8 2.3 Utility: Steering Score

SAE steering injects the SAE decoder atom  $v_f$  (the f-th column of the column-normalized decoder  $W_{\text{dec}}[f]$ ) into the residual stream at a target layer to push the hidden state x along a chosen feature direction [Durmus et al., 2024]. Given a feature index f, a steering factor  $\alpha$ , and a per-sample scale  $m_f$  (e.g., the feature's maximum activation), the intervention is

$$x^{\text{steer}} = x + (\alpha m_f) \cdot v_f. \tag{1}$$

Through the above formula (1), we can use SAE features for steering to achieve the output of controlling LLM. AXBENCH [Wu et al., 2025] measures causal control by steering internal representations during generation and asking an LLM judge to rate three aspects, each on a discrete scale  $\{0,1,2\}$ : Concept (C), Instruction (I), and Fluency (F). The overall Steering Score is the harmonic mean:

Steering Score = 
$$\operatorname{HM}(C, I, F) = \frac{3}{\frac{1}{C} + \frac{1}{I} + \frac{1}{F}} \in [0, 2].$$

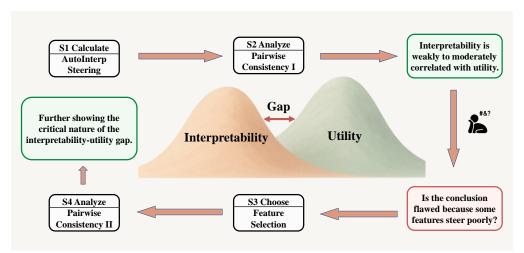


Figure 2: Overview of our pairwise-controlled workflow linking SAE interpretability with steering utility. (S1) Compute interpretability score and steering score for each SAE. (S2) Pairwise analysis across SAEs and get an insight (the top-right green box), revealing an interpretability—utility gap. The red box (lower right) is our further inference based on the above green box and previous studies [Wu et al., 2025]. (S3) Use  $\Delta$  Token Confidence to select higher-utility features. (S4) Compute steering gains after selection per SAE, then do the pairwise analysis between steering gains and interpretability. The green box in the middle left is our final conclusion.

Following AXBENCH, for each concept we sample instructions (e.g., 10 from Alpaca-Eval [Dubois et al., 2023]), generate continuations under different steering factors, pick the best factor on one split, and evaluate the held-out split with the judge to obtain the final utility score averaged across prompts [Gu et al., 2025]. The complete scoring procedure is detailed in Appendix C.

## 3 Can SAE Interpretability Indicate Steering Performance?

#### 3.1 Experimental Setup

112

126

Dataset. For each trained SAE, we score 1,000 latents with *LLM-as-judge* [Paulo et al., 2025] and randomly sample 100 to form that SAE's CONCEPT100 (see Appendix F). For steering, we sample 10 Alpaca-Eval instructions, allow up to 128 generated tokens, and test 6 steering factors; the 10 instructions are split 5/5 for factor selection vs. held-out evaluation.

Model. We evaluate three open LLMs: Gemma-2-2B [Team et al., 2024], Qwen-2.5-3B [Yang et al., 2024], and Gemma-2-9B [Team et al., 2024]. SAEs are trained on residual-stream activations at a fixed mid-layer for each model: Layer 12 for Gemma-2-2B, Layer 17 for Qwen-2.5-3B, and Layer 20 for Gemma-2-9B—and steering is applied to the corresponding layer.

SAE with different architectures We train 90 SAEs covering a range of architectures and sparsity. All SAEs use a latent dictionary width of 16k. We instantiate five variants: BatchTopK [Bussmann et al., 2024], Gated [Rajamanoharan et al., 2024a], JumpReLU [Rajamanoharan et al., 2024b], ReLU [Team, 2024], TopK [Gao et al., 2024] and sweep six target sparsity levels with approximate per-token activations  $L_0 \approx 50, 80, 160, 320, 520, 820$ . Further details are provided in Appendix A.

## 3.2 Pairwise Rank Consistency between Interpretability and Utility

We test whether higher interpretability of SAE is predictive of higher steering performance across a set of trained SAEs attached to a fixed LM. For each SAE  $\theta$  in a pool  $\Theta$ , we record a pair  $(\mu(\theta), g(\theta)) \in \mathbb{R}^2$ , where  $\mu$  is the SAE-level *Interpretability Score* and g is an aggregated *Steering Score* over a standardized evaluation suite.

Given two SAEs  $\theta_i, \theta_i \in \Theta$ , define the concordance indicator

$$v_{ij} = \operatorname{sign}(\mu(\theta_i) - \mu(\theta_j)) \cdot \operatorname{sign}(g(\theta_i) - g(\theta_j)) \in \{-1, 0, +1\}. \tag{2}$$

Kendall's tie-corrected rank coefficient  $\tau_b$  [KENDALL, 1938] summarizes agreement over unordered pairs and reduces to average concordance when there are no ties:

$$\tau_b = \frac{1}{\binom{|\Theta|}{2}} \sum_{i < j} v_{ij} \in [-1, 1]. \tag{3}$$

In this study, we instantiate  $\mu$  with the *Interpretability Score* and g with the *Steering Score*, then compute  $\tau$  for three model–layer settings (Gemma-2-2B, Qwen-2.5-3B, Gemma-2-9B). Each setting includes 30 SAEs spanning architectures and sparsity to ensure sufficient pair coverage.

## 3.3 Granulated Kendall's Coefficient to Control Confounders

between groups to obtain the statistic at the axis level:

and utility. To obtain an axis-controlled assessment, we factor the SAE design space into orthogonal axes and evaluate rank consistency while varying one axis at a time and holding the others fixed.

We define three conditioning axes: (A) Architecture — fix architecture (and layer), vary sparsity;

(B) Sparsity — compare architectures at matched sparsity ranks; (C) Model — fix the base model, compare all SAEs within it. For axis i, partition  $\Theta$  into groups  $\mathcal{G}_i$  that are matched on all axes except i. Within each group  $G \in \mathcal{G}_i$ , compute Kendall's coefficient in  $\{(\mu(\theta), g(\theta)) : \theta \in G\}$ , and average

Global rank agreement can be confounded by hyperparameters that jointly influence interpretability

$$\psi_i = \frac{1}{|\mathcal{G}_i|} \sum_{G \in \mathcal{G}_i} \tau(\{(\mu(\theta), g(\theta)) : \theta \in G\}). \tag{4}$$

146 Aggregate the axis-level outcomes by

$$\Psi = \frac{1}{n} \sum_{i=1}^{n} \psi_i, \tag{5}$$

where n is the number of axes. Each  $\psi_i$  captures rank consistency conditioned on axis i (varying only that axis while matching the others), and  $\Psi$  aggregates these into a single axis-controlled measure. This construction mitigates cross-axis trends—e.g., architecture, sparsity, or model-driven shifts that can obscure local relationships between interpretability and utility.

We report the per-axis statistics  $\psi_i$  together with the aggregate  $\Psi$  for the same model settings as in section 3.2, providing both axis-specific and aggregated assessments.

## 3.4 Pairwise Analysis Results

In this section, we assess whether higher SAE interpretability predicts stronger steering by computing Kendall's  $\tau_b$  between the *Interpretability Score*  $\mu(\theta)$  and the aggregated *Steering Score*  $g(\theta)$  over a pooled set of SAEs attached to a fixed LLM.

To control confounders and localize effects, we apply the axis-conditioned procedure defined in section 3.3. For each axis, we form matched groups, compute within-group  $\tau_b$ , average to obtain a per-axis summary, and aggregate these summaries into an overall axis-controlled coefficient.

Table 1 shows that across SAEs, higher interpretability tends to be modestly associated with better steering on average, pointing to a consistent but limited impact. The pooled Kendall's  $\tau_b \approx 0.30$  is positive, and the axis-controlled aggregate remains positive ( $\Psi \approx 0.25$ ), indicating that more interpretable features generally translate into better steering utility across designs and models.

The strength of the link between interpretability and utility depends on SAE architecture, sparsity, and the base model. By architecture, the association is positive on average ( $\Psi_A \approx 0.26$ ), with ReLU-like variants reinforcing the trend and Gated weakening it. By sparsity, alignment is strongest when the SAE is more sparse and weakens—sometimes reversing—as the number of active features increases. By model, the underlying LM shapes the effect, with the signal clearest in Qwen-2.5-3B and weaker in Gemma-2-2B, while the model-wise summary remains positive ( $\Psi_C \approx 0.33$ ).

**Key Observation 1:** Interpretability shows a relatively weak positive correlation with steering performance, highlighting a notable gap between interpretability and utility across SAEs.

164

165

138

145

Table 1: Pairwise Analysis Between Interpretability Score  $\mu(\theta)$  and Steering Score  $g(\theta)$ . We report Kendall's  $\tau_b$  overall and by axis-controlled measures  $\Psi_A$  (Architecture),  $\Psi_B$  (Matched Sparsity), and  $\Psi_C$ (Model). n = number of SAEs; Pairs = number of pairwise comparisons; p = permutation p-value; 95% CI = confidence interval (overall uses BCa bootstrap; subgroups use permutation-based CIs).

Axis	SAEs	n	Pairs	$ au_b$	p	95% CI
Overall	All SAEs	90	4005	0.2979	_	[0.1590, 0.4191]
	$\Psi_{ m A}$	= 0.257	$5 \text{ (SE} \approx 0.1)$	163, 95% bo	ot CI [0.0222	2, 0.3961])
	BatchTopK	18	153	0.3203	0.0712	[-0.3464, 0.3464]
	Gated	18	153	-0.2026	0.2577	[-0.3464, 0.3333]
$\Psi_{\mathbf{A}}$ : Architecture	JumpReLU	18	153	0.4248	0.0160	[-0.3333, 0.3333]
	ReLU	18	153	0.3595	0.0392	[-0.3333, 0.3333]
	TopK	18	153	0.3856	0.0272	[-0.3333, 0.3333]
	$\Psi_{\mathrm{B}}$ =	= 0.1651	$(SE \approx 0.1)$	112, 95% boo	t CI [-0.028	36, 0.3587])
	$L_0 \approx 50$	15	105	0.5429	0.0034	[-0.3714, 0.3714]
	$L_0 \approx 80$	15	105	0.3524	0.0740	[-0.3714, 0.3714]
$\Psi_{\rm B}$ : Sparsity	$L_0 \approx 160$	15	105	0.1810	0.3821	[-0.3905, 0.3714]
ΨB. Sparsity	$L_0 \approx 320$	15	105	0.1810	0.3673	[-0.3714, 0.3714]
	$L_0 \approx 520$	15	105	-0.2190	0.2837	[-0.3905, 0.3714]
	$L_0 \approx 820$	15	105	-0.0476	0.8484	[-0.3905, 0.3714]
	$\Psi_{\mathrm{C}}$	= 0.327	$2 \text{ (SE} \approx 0.0$	0698, 95% bo	ot CI [0.2184	4, 0.4575])
	Gemma-2-2B	30	435	0.2184	0.0980	[-0.2598,  0.2552]
$\Psi_{\mathbf{C}}$ : Model	Qwen-2.5-3B	30	435	0.4575	0.0008	[-0.2506, 0.2506]
	Gemma-2-9B	30	435	0.3057	0.0166	[-0.2506, 0.2461]
		Л	$\Psi = (\Psi_{\rm A} +$	$\Psi_{\rm B} + \Psi_{\rm C})/3$	8 = 0.2499	

## From Interpretability to Utility: Which SAE Features Actually Steer?

In Sec. 3.4, We find that SAE interpretability is a relatively weak prior for steering utility. Prior 172 work [Arad et al., 2025] shows many features lack steerability and we speculate that this factor may render the previous conclusion inaccurate. Therefore, we introduce a metric to identify steeringeffective features. Metrics derived from a model's internal token distributions can assess reasoning quality [Kang et al., 2025]. In particular, token entropy offers a unified view: high entropy highlights critical decision points [Fu et al., 2025, Wang et al., 2025c]. We apply this idea to SAE steering. 177

#### 4.1 Feature Selection via $\Delta$ Token Confidence 178

173

174

175

176

We start from the model's next-token distribution. Given logits  $z \in \mathbb{R}^V$  and  $p = \operatorname{softmax}(z)$  over a vocabulary of size V, the token entropy is

$$H(p) = -\sum_{j=1}^{V} p_j \log p_j,$$
 (6)

Entropy summarizes dispersion over the vocabulary: smaller values reflect a sharper, more concen-181 trated prediction, while larger values indicate greater uncertainty at a given position. 182

To focus on the head of the distribution that matters most for sampling, we use token confidence. 183 Let  $\mathcal{I}_k(p) \subseteq \{1,\ldots,V\}$  denote the indices of the k largest probabilities in p. The top-k token 184 confidence is the negative average log-probability over these entries: 185

$$C_k(p) = -\frac{1}{k} \sum_{j \in \mathcal{I}_k(p)} \log p_j. \tag{7}$$

Lower  $C_k$  implies higher confidence, while higher  $C_k$  implies a flatter top-k distribution. Unlike 186 entropy,  $C_k$  directly captures the sharpness of the outcomes that drive next-token behavior. 187

We turn confidence into a feature-level selector via a single-feature SAE intervention. Consider an SAE feature f at layer  $\ell$ . We amplify only the coefficient of f by a factor  $\alpha > 0$  in the SAE reconstruction, leaving the base model and all other features unchanged. Denote the baseline next-

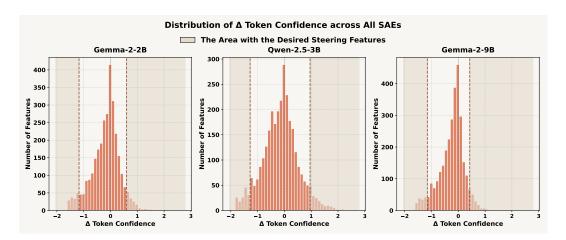


Figure 3: **Distribution of per-feature**  $\Delta$  *Token Confidence* across all SAEs. Panels show histograms for Gemma-2-2B, Qwen-2.5-3B, and Gemma-2-9B; the x-axis is  $\Delta C_k$  (negative values indicate increased confidence, positive values decreased confidence) and the y-axis is the number of SAE features. The shaded area marks the high-magnitude tails from which candidate steering features are selected, while the central mass near 0 indicates features with little distributional impact.

token distribution by  $p^{\mathrm{base}}$  and the intervened distribution by  $p^{\mathrm{int}}_{f,\ell,\alpha}$ .  $\Delta$  Token Confidence is

$$\Delta C_k(f;\ell,\alpha) = C_k(p_{f,\ell,\alpha}^{\text{int}}) - C_k(p^{\text{base}}).$$
(8)

Negative values  $\Delta C_k < 0$  mean that amplifying f sharpens the top-k distribution, while positive values indicate greater dispersion. We compute this using one baseline and one intervened forward pass via an SAE hook. Implementation details and hyperparameters are provided in Appendix D.

We select features with the largest absolute change in *token confidence* under single-feature SAE interventions, i.e., maximal  $|\Delta C_k|$  (see Figure 3). For each feature, we compute  $\Delta C_k$ , rank by  $|\Delta C_k|$ , form tiers, evaluate subsets for steering, and keep the best per SAE.

## 4.2 Steering Performance Results After Feature Selection

Arad et al. [2025] has shown that SAE steering works well if features are chosen by their causal impact on model outputs, introducing the *output score* as a metric to identify output-aligned features. Following this insight, we evaluate our  $\Delta$  *token confidence* selection on three base LLMs (Gemma-2-2B, Qwen-2.5-3B, Gemma-2-9B) using the CONCEPT100 (see details in 3.1). The experiments on steering performance improvement of each SAE can be referred to Appendix E.

Table 2 shows that our selection yields consistent gains across all models, outperforming the vanilla SAE baseline by large margins, and also improving over an output-score—based selector. These gains indicate that ranking and filtering by the magnitude of distributional change captured by  $\Delta C_k$  reliably isolates features with the strongest steering utility.

Furthermore, we conducted a comparative analysis of SAEs of different architectures on three models. For fair comparison, the two feature selection methods use the same subset size. Figure 4 compares *steer*-

Table 2: Steering score after feature selection compared with SAE-based steering. Columns report scores (higher is better) for Gemma-2-2B, Qwen-2.5-3B, and Gemma-2-9B. Rows: 'SAE-based' uses all SAE features without selection [Wu et al., 2025]; '+Output' selects features using  $S_{\rm out}$  [Arad et al., 2025]; "+ $\Delta C_k$  (Ours)" selects by the  $\Delta$  Token Confidence. Boldface indicates the best method per model.

Method	Gemma-2-2B	<b>Qwen-2.5-3B</b>	Gemma-2-9B
SAE-based	0.133	0.171	0.142
+Output	0.233	0.292	0.255
$+\Delta C_k(\text{Ours})$	0.328	0.399	0.289

 $ing\ scores$  across SAE architectures and selection methods. In all three models, selecting features by  $\Delta\ Token\ Confidence$  consistently outperforms both the no-selection SAE baseline and the output-score selector across architectures.

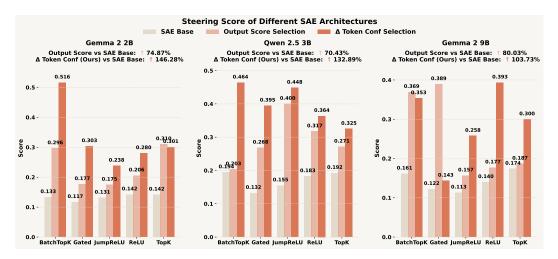


Figure 4: Comparison of different SAE steering methods with five SAE architecture across three LLMs. Panels correspond to Gemma-2-2B, Qwen-2.5-3B, and Gemma-2-9B. The horizontal axis groups SAE architectures (BatchTopK, Gated, JumpReLU, ReLU, TopK), and the vertical axis reports the *steering score*. Bars show three conditions: SAE Base (no feature selection), Output Score Selection, and  $\Delta$  *Token Confidence* Selection (ours). Panel annotations summarize the average lift of each selection method relative to the SAE-based steering.

On average, our method improves steering performance by **52.52**% over the strongest competing baseline. The BatchTopK architecture is the one that has the most stable and significant improvement in steering capabilities on models of different sizes among the five SAE architectures.

**Key Observation 2:**  $\Delta$  *Token Confidence* reliably selects high-utility SAE features across models. Among SAE architectures, BatchTopK achieves the most stable and sizable *steering gains*.

## 4.3 Pairwise Analysis Between Interpretability and Steering Gain

Building on the above *steering gains*, we further examine whether SAE interpretability can serve as a prior for *steering gain*, the extent to which a trained SAE benefits from feature selection. We quantify this relationship by computing Kendall's  $\tau_b$  between the *Interpretability Score*  $\mu(\theta)$  of SAEs and the *Steering Gain*  $L(\theta)$ , defined as the percentage lift of the selected-steering score over the same SAE's base. As in section 3.4, we report both pooled coefficients and axis-conditioned summaries that control for design and model factors (architecture, sparsity and model).

Overall, table 3 indicates that interpretability is not a reliable prior for *steering gain* after selection: the pooled association is small and slightly negative ( $\tau_b \approx -0.069$ ), and the axis-controlled aggregate is likewise near zero and negative ( $\Psi \approx -0.057$ ). Estimates cluster near zero across design axes, being slightly negative within architecture, sparsity, and model, and effectively null at matched-sparsity slots, indicating no consistent link between higher interpretability and larger gains.

**Key Observation 3:** Surprisingly, the interpretability–utility gap widens when we focus on SAE features that deliver substantial *steering gains*.

## 5 Related Work

## 5.1 Representation-Based Steering

Activation-based steering arose as a lightweight alternative to fine-tuning, enabling on-the-fly control of LLM behavior without retraining [Giulianelli et al., 2018, Vig et al., 2020, Geiger et al., 2021, 2025]. The core idea is to inject carefully chosen directions into hidden states, typically in the residual stream, scaling interventions by a gain and selecting layers for maximal effect [Zou et al., 2025, Rimsky et al., 2024, van der Weij et al., 2024]. It has been applied to safety and moderation,

Table 3: Pairwise Analysis Between Interpretability  $\mu(\theta)$  and Steering Gain  $L(\theta)$ . We report Kendall's  $\tau_b$  overall and under axis-controlled summaries  $\Psi_A$  (Architecture),  $\Psi_B$  (Matched Sparsity), and  $\Psi_C$  (Model). Columns: n = number of SAEs; Pairs = number of pairwise comparisons; p = permutation p-value; 95% CI = confidence interval (overall uses BCa bootstrap; subgroups use permutation-based CIs).

Axis	SAEs	n	Pairs	$ au_b$	p	95% CI
Overall	All SAEs	90	4005	-0.0692	_	[-0.2019, 0.0666]
	$\Psi_{ m A} =$	-0.0719	9 (SE $\approx 0$ .	0781, 95% b	oot CI $[-0.$	2078, 0.0614])
	BatchTopK	18	153	-0.2288	0.2004	[-0.3333, 0.3464]
	Gated	18	153	0.0327	0.8792	[-0.3464, 0.3333]
$\Psi_{\mathbf{A}}$ : Architecture	JumpReLU	18	153	-0.0065	1.0000	[-0.3464, 0.3203]
	ReLU	18	153	-0.2810	0.1096	[-0.3333, 0.3464]
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	0.4985	[-0.3464, 0.3337]				
$\Psi_{\rm B} =$ <b>0.0127</b> (SE $\approx 0.0457, 95\%$ boot CI $[-0.0762, 0.0889]$ )						
	$L_0 \approx 50$	15	105	0.0476	0.8466	[-0.3714, 0.3714]
	$L_0 \approx 80$	15	105	0.1619	0.4437	[-0.3905, 0.3714]
II. · Sparcity	$L_0 \approx 160$	15	105	0.0095	1.0000	[-0.3714, 0.3714]
ΨB. Sparsity	$L_0 \approx 320$	15	105	0.0476	0.8452	[-0.3714, 0.3714]
	$L_0 \approx 520$	15	105	-0.1810	0.3797	[-0.3714, 0.3714]
	$L_0 \approx 820$	15	105	-0.0095	1.0000	[-0.3714, 0.3905]
	$\Psi_{\mathrm{C}} = 0$	-0.1111	$(SE \approx 0.0)$	314, 95% bo	oot CI [-0.1	[448, -0.0483])
	Gemma-2-2B	30	435	-0.1402	0.2911	[-0.2552, 0.2598]
$\Psi_{\mathbf{C}}$ : Model	Qwen-2.5-3B	30	435	-0.1448	0.2781	[-0.2507, 0.2460]
	Gemma-2-9B	30	435	-0.0483	0.7157	[-0.2506, 0.2552]
		Ŋ	$\Psi = (\Psi_{ m A} +$	$-\Psi_{\rm B}+\Psi_{\rm C}$	/3 = -0.05	568

persona and sentiment control, and instruction adherence, promising low-latency deployment-time adjustment but facing polysemantic entanglement and brittleness that motivate standardized evaluation [Chen et al., 2025, Liu et al., 2024]. However, this approach injects polysemantic activations at intervention time, yielding coarse-grained effects for output control [Bricken et al., 2023]. Our work is related to activation-level interventions, but differs by grounding directions in sparse, interpretable SAE features and applying utility-oriented feature selection to mitigate these failure modes.

#### 5.2 SAE-Based Steering

Sparse Autoencoders (SAEs) decompose activations into sparse, human-readable features to mitigate polysemanticity and expose concept-level structure [Bricken et al., 2023, Templeton et al., 2024, Gao et al., 2024]. For steering, practitioners use decoder atoms as directions and add scaled injections at chosen layers, with architecture and sparsity choices trading reconstruction for feature granularity [Zhao et al., 2025, Wang et al., 2025d, Ferrando et al., 2025]. SAE-based steering enables targeted safety control, style modulation, and instruction emphasis, yet the utility of individual features varies widely [Chalnev et al., 2024, Mayne et al., 2024]. While the connection between SAE interpretability and steering utility remains unclear, and our goal is to build a principled bridge between them. To this end, we conduct a large-scale experiments across multiple model sizes and SAE architectures, demonstrating the critical nature of the interpretability-utility gap.

## Conclusion and Discussion

In summary, SAE interpretability shows relatively weak positive association with steering utility across 90 SAEs ( $\tau_b \approx 0.298$ ), revealing a clear interpretability–utility gap. Selecting features with  $\Delta$  Token Confidence yields substantial gains (average +52.52% over the strongest existing baseline). Surprisingly, when analyzing steering gains after selection, the correlation with interpretability collapses toward zero and can even turn negative for the highest-utility features, further underscoring this gap. This gap points to a key direction: develop task-general utility indicators that reliably predict steerability across models, or design training objectives that directly optimize controllability under sparsity so features are utility-calibrated without heavy post-hoc selection. Our work provides valuable insight for the further development of SAEs as interpretable tools.

#### 272 Reproducibility Statement

We aim to facilitate full reproduction of our results. All model code, training and evaluation scripts, 273 configuration files, and experiment logs are released at an anonymous repository as part of the sup-274 plementary materials: https://anonymous.4open.science/r/SAE4Steer. Training architectures, hy-275 perparameters, sparsity schedules, and optimization details are specified in the main text and Appendix A.2 (see also the per-family settings in Appendix A). The datasets used are openly licensed: all SAEs are trained on The Common Pile v0.1 [Kandpal et al., 2025] as described in Appendix F; our evaluation concepts (CONCEPT100) and their automatic generation pipeline are documented in Appendix B and Appendix F. The complete procedures for automated interpretability scoring 280 (SAEBENCH) and steering utility (AXBENCH), including sampling, judging protocols, and scoring 281 functions, are detailed in Appendix B and Appendix C, with the  $\Delta$  Token Confidence selector defined 282 in Appendix D and the post-selection results summarized in Appendix E. Hardware, runtime, and 283 memory footprints for both SAEBENCH and AXBENCH are reported in Appendices B.3 and C.2. 284 Together, these materials, along with seed-controlled configuration files and exact command-line 285 286 invocations provided in the anonymous repository, are intended to enable independent researchers to replicate and extend our findings. 287

## References

- Adam Karvonen, Can Rager, Johnny Lin, Curt Tigges, Joseph Isaac Bloom, David Chanin, YeuTong Lau, Eoin Farrell, Callum Stuart McDougall, Kola Ayonrinde, Demian Till, Matthew Wearden, Arthur Conmy, Samuel Marks, and Neel Nanda. SAEBench: A comprehensive benchmark
  for sparse autoencoders in language model interpretability. In *Forty-second International Confer-*ence on Machine Learning, 2025. URL https://openreview.net/forum?id=qrU3yNfX0d.
- Zhengxuan Wu, Aryaman Arora, Atticus Geiger, Zheng Wang, Jing Huang, Dan Jurafsky, Christopher D Manning, and Christopher Potts. Axbench: Steering LLMs? even simple baselines outperform sparse autoencoders. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=K2CckZjNy0.
- Dana Arad, Aaron Mueller, and Yonatan Belinkov. Saes are good for steering if you select the right features, 2025. URL https://arxiv.org/abs/2505.20063.
- Aounon Kumar, Chirag Agarwal, Suraj Srinivas, Aaron Jiaxun Li, Soheil Feizi, and Himabindu Lakkaraju. Certifying LLM safety against adversarial prompting. In *Proceedings of the Conference on Language Modeling (COLM)*, 2024. URL https://arxiv.org/abs/2309.02705.
- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang
  Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment
  of Ilm via a human-preference dataset. In *Advances in Neural Information Processing Sys-*tems (NeurIPS), 2023. URL http://papers.nips.cc/paper\_files/paper/2023/hash/
  4dbb61cb68671edc4ca3712d70083b9f-Abstract-Datasets\_and\_Benchmarks.html.
- Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madian Khabsa. Llama guard: Llmbased input-output safeguard for human-ai conversations, 2023. URL https://arxiv.org/abs/2312.06674.
- Yihan Chen, Benfeng Xu, Quan Wang, Yi Liu, and Zhendong Mao. Benchmarking large language models on controllable generation under diversified instructions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17808–17816, 2024.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9.
- Xu Wang, Yan Hu, Wenyu Du, Reynold Cheng, Benyou Wang, and Difan Zou. Towards understanding fine-tuning mechanisms of LLMs via circuit analysis. In *Forty-second International Conference on Machine Learning*, 2025a. URL https://openreview.net/forum?id=45EIiFd60a.

- Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. Steering Language Models With Activation Engineering. *arXiv e-prints*, art. arXiv:2308.10248, August 2023. doi: 10.48550/arXiv.2308.10248.
- Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. Steering language models with activation engineering, 2024. URL https://arxiv.org/abs/2308.10248.
- Tianlong Wang, Xianfeng Jiao, Yinghao Zhu, Zhongzhi Chen, Yifan He, Xu Chu, Junyi Gao, Yasha Wang, and Liantao Ma. Adaptive activation steering: A tuning-free llm truthfulness improvement method for diverse hallucinations categories. In *Proceedings of the ACM on Web Conference* 2025, pages 2562–2578, 2025b.
- Alessandro Stolfo, Vidhisha Balachandran, Safoora Yousefi, Eric Horvitz, and Besmira Nushi. Improving instruction-following in language models through activation steering. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=wozhdnRCtw.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. https://transformercircuits.pub/2023/monosemantic-features/index.html.
- Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam
  Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner,
  Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees,
  Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*,
  2024. URL https://transformer-circuits.pub/2024/scaling-monosemanticity/
  index.html.
- Anish Mudide, Joshua Engels, Eric J. Michaud, Max Tegmark, and Christian Schroeder de Witt.

  Efficient dictionary learning with switch sparse autoencoders, 2025. URL https://arxiv.
  org/abs/2410.08201.
- Javier Ferrando, Oscar Balcells Obeso, Senthooran Rajamanoharan, and Neel Nanda. Do i know this entity? knowledge awareness and hallucinations in language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=WCRQFlji2q.
- Sviatoslav Chalnev, Matthew Siu, and Arthur Conmy. Improving steering vectors by targeting sparse autoencoder features, 2024. URL https://arxiv.org/abs/2411.02193.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models, 2023. URL https://arxiv.org/ abs/2309.08600.
- Aleksandar Makelov. Sparse autoencoders match supervised features for model steering on the IOI task. In *ICML 2024 Workshop on Mechanistic Interpretability*, 2024. URL https://openreview.net/forum?id=JdrVuEQih5.
- Kyle O'Brien, David Majercak, Xavier Fernandes, Richard G Edgar, Blake Bullwinkel, Jingya Chen,
   Harsha Nori, Dean Carignan, Eric Horvitz, and Forough Poursabzi-Sangdeh. Steering language
   model refusal with sparse autoencoders. In *ICML 2025 Workshop on Reliable and Responsible Foundation Models*.
- Gonçalo Santos Paulo, Alex Troy Mallen, Caden Juang, and Nora Belrose. Automatically interpreting millions of features in large language models. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=EemtbhJOXc.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu-372 patiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Fer-373 ret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Char-374 line Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, 375 Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, 376 Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchi-377 son, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, 378 Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, 379 Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Wein-380 berger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, 381 Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, 382 Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen 383 Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha 384 Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kar-386 tikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, 387 Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, 388 Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel 389 Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, 390 Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moyni-391 han, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, 392 Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil 393 Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culli-394 ton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, 395 Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Cogan, Sarah Perrin, 396 Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ron-397 strom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee 398 Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei 399 Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan 400 Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli 401 Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dra-402 gan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Fara-403 bet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, 404 Robert Dadashi, and Alek Andreev. Gemma 2: Improving open language models at a practical 405 size, 2024. URL https://arxiv.org/abs/2408.00118. 406

Qwen An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan 407 Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, 408 Jianwei Zhang, Jianxin Yang, Jiaxin Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, 409 Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji 410 Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yi-Chao 411 Zhang, Yunyang Wan, Yuqi Liu, Zeyu Cui, Zhenru Zhang, Zihan Qiu, Shanghaoran Quan, and 412 Zekun Wang. Qwen2.5 technical report. ArXiv, abs/2412.15115, 2024. URL https://api. 413 semanticscholar.org/CorpusID:274859421. 414

Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=SJgIPJBFvH.

Yunzhe Hu, Difan Zou, and Dong Xu. An in-depth investigation of sparse rate reduction in transformer-like models. *Advances in Neural Information Processing Systems*, 37:116815–116837, 2024.

Yichao Fu, Xuewei Wang, Yuandong Tian, and Jiawei Zhao. Deep think with confidence, 2025. URL https://arxiv.org/abs/2508.15260.

Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen,
Jianxin Yang, Zhenru Zhang, Yuqiong Liu, An Yang, Andrew Zhao, Yang Yue, Shiji Song, Bowen
Yu, Gao Huang, and Junyang Lin. Beyond the 80/20 rule: High-entropy minority tokens drive
effective reinforcement learning for llm reasoning, 2025c. URL https://arxiv.org/abs/
2506.01939.

- Patrick Leask, Bart Bussmann, Michael T Pearce, Joseph Isaac Bloom, Curt Tigges, Noura Al Moubayed, Lee Sharkey, and Neel Nanda. Sparse autoencoders do not find canonical units of analysis. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=9ca9eHNrdH.
- Adam Karvonen, Benjamin Wright, Can Rager, Rico Angell, Jannik Brinkmann, Logan Smith,
  Claudio Mayrink Verdun, David Bau, and Samuel Marks. Measuring progress in dictionary learning for language model interpretability with board game models. In A. Globerson,
  L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, Advances in
  Neural Information Processing Systems, volume 37, pages 83091–83118. Curran Associates,
  Inc., 2024. URL https://proceedings.neurips.cc/paper\_files/paper/2024/file/
  9736acf007760cc2b47948ae3cf06274-Paper-Conference.pdf.
- Esin Durmus, Alex Tamkin, Jack Clark, Jerry Wei, Jonathan Marcus, Joshua Batson, Kunal Handa, Liane Lovitt, Meg Tong, Miles McCain, Oliver Rausch, Saffron Huang, Sam Bowman, Stuart Ritchie, Tom Henighan, and Deep Ganguli. Evaluating feature steering: A case study in mitigating social biases, 2024. URL https://anthropic.com/research/ evaluating-feature-steering.
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos
   Guestrin, Percy Liang, and Tatsunori B Hashimoto. Alpacafarm: a simulation framework for
   methods that learn from human feedback. In *Proceedings of the 37th International Conference* on Neural Information Processing Systems, pages 30039–30069, 2023.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan
   Shen, Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel
   Ni, and Jian Guo. A survey on llm-as-a-judge, 2025. URL https://arxiv.org/abs/2411.
   15594.
- Bart Bussmann, Patrick Leask, and Neel Nanda. Batchtopk: A simple improvement for topksaes.
   In AI Alignment Forum, page 17, 2024.
- Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Tom Lieberum, Vikrant Varma, János
   Kramár, Rohin Shah, and Neel Nanda. Improving dictionary learning with gated sparse autoen coders, 2024a. URL https://arxiv.org/abs/2404.16014.
- Senthooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse autoencoders, 2024b. URL https://arxiv.org/abs/2407.14435.
- Anthropic Interpretability Team. Training sparse autoencoders. https://
  transformer-circuits.pub/2024/april-update/index.html#training-saes, 2024.
  Accessed: 2025-01-20.
- Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever,
  Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders, 2024. URL https://arxiv.org/abs/2406.04093.
- M. G. KENDALL. A new measure of rank correlation. *Biometrika*, 30(1-2):81-93, 06 1938. ISSN 0006-3444. doi: 10.1093/biomet/30.1-2.81. URL https://doi.org/10.1093/biomet/30.1-2.81.
- Zhewei Kang, Xuandong Zhao, and Dawn Song. Scalable best-of-n selection for large language
   models via self-certainty, 2025. URL https://arxiv.org/abs/2502.18581.
- Mario Giulianelli, Jack Harding, Florian Mohnert, Dieuwke Hupkes, and Willem Zuidema. Under the hood: Using diagnostic classifiers to investigate and improve how language models track agreement information. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 240–248, 2018.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. Investigating gender bias in language models using causal mediation analysis. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in*

- Neural Information Processing Systems, volume 33, pages 12388-12401. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper\_files/paper/2020/file/92650b2e92217715fe312e6fa7b90d82-Paper.pdf.
- Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. Causal abstractions of neural networks. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, Advances in Neural Information Processing Systems, volume 34, pages 9574–9586. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper\_files/paper/2021/file/4f5c422f4d49a5a807eda27434231040-Paper.pdf.
- Atticus Geiger, Duligur Ibeling, Amir Zur, Maheep Chaudhary, Sonakshi Chauhan, Jing Huang,
  Aryaman Arora, Zhengxuan Wu, Noah Goodman, Christopher Potts, and Thomas Icard. Causal
  abstraction: A theoretical foundation for mechanistic interpretability. *Journal of Machine Learn-ing Research*, 26(83):1–64, 2025. URL http://jmlr.org/papers/v26/23-0058.html.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan Hendrycks. Representation engineering: A top-down approach to ai transparency, 2025. URL https://arxiv.org/abs/2310.01405.
- Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. Steering llama 2 via contrastive activation addition. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15504–15522, 2024.
- Teun van der Weij, Massimo Poesio, and Nandi Schoots. Extending activation steering to broad skills and multiple behaviours, 2024. URL https://arxiv.org/abs/2403.05767.
- Runjin Chen, Andy Arditi, Henry Sleight, Owain Evans, and Jack Lindsey. Persona vectors: Monitoring and controlling character traits in language models, 2025. URL https://arxiv.org/ abs/2507.21509.
- Sheng Liu, Haotian Ye, Lei Xing, and James Y. Zou. In-context vectors: Making in context learning more effective and controllable through latent space steering. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 32287–32307. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/liu24bx.html.
- Yu Zhao, Alessio Devoto, Giwon Hong, Xiaotang Du, Aryo Pradipta Gema, Hongru Wang, Xuanli
  He, Kam-Fai Wong, and Pasquale Minervini. Steering knowledge selection behaviours in LLMs
  via SAE-based representation engineering. In Luis Chiruzzo, Alan Ritter, and Lu Wang, editors,

  Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association
  for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages
  514
  515 URL https://aclanthology.org/2025.naacl-long.264/.
- Xu Wang, Zihao Li, Benyou Wang, Yan Hu, and Difan Zou. Model unlearning via sparse autoencoder subspace guided projections. In *ICML 2025 Workshop on Machine Unlearning for Generative AI*, 2025d. URL https://openreview.net/forum?id=MIlqM98o9I.
- Harry Mayne, Yushi Yang, and Adam Mahdi. Can sparse autoencoders be used to decompose and interpret steering vectors?, 2024. URL https://arxiv.org/abs/2411.08790.
- Nikhil Kandpal, Brian Lester, Colin Raffel, Sebastian Majstorovic, Stella Biderman, Baber Abbasi, Luca Soldaini, Enrico Shippole, A. Feder Cooper, Aviya Skowron, Shayne Longpre, Lintang Sutawika, Alon Albalak, Zhenlin Xu, Guilherme Penedo, Loubna Ben Allal, Elie Bakouch, John David Pressman, Honglu Fan, Dashiell Stander, Guangyu Song, Aaron Gokaslan, John Kirchenbauer, Tom Goldstein, Brian R. Bartoldson, Bhavya Kailkhura, and Tyler Murray. The Common Pile v0.1: An 8TB Dataset of Public Domain and Openly Licensed Text. arXiv preprint arXiv:2506.05209, 2025.

## 28 LLM Usage

537

- 529 In preparing this paper, large language models (LLMs) were used as an assistive tool for minor
- sao language polishing and stylistic improvements. All technical contributions, results, and conclusions
- are solely the work of the authors.

## 532 A SAE Architectures and Training Details

- We train 90 SAEs (30 per base model) across five architectures and six target sparsity levels. Unless
- stated otherwise, the dictionary width is 16K codes (F=16,384), SAEs are attached to the residual
- stream at the layer described in the main text, and decoder columns are  $\ell_2$ -normalized. All models
- are trained on The Common Pile v0.1 [Kandpal et al., 2025].

#### A.1 Architectures and Parameterization

We list the five SAE families with their named parameters (as implemented) and the corresponding shapes. The last column records architecture-specific thresholding/gating fields when present.

Shapes assume residual dimension d=2304 and dictionary width F=16,384.

Architectures	$W_{ m enc}$	$b_{ m enc}$	$W_{ m dec}$	$b_{ m dec}$	Threshold / Extras
ReLU	encoder.weight: shape (16,384, 2,304)	encoder.bias: shape (16,384)	decoder.weight: shape (2,304, 16,384)		_
Gated	encoder.weight: shape (16,384, 2,304)	gate_bias: shape (16,384)	decoder.weight: shape (2,304, 16,384)		r_mag: shape (16,384); mag_bias: shape (16,384)
TopK	encoder.weight: shape (16,384, 2,304)	encoder.bias: shape (16,384)	decoder.weight: shape (2,304, 16,384)		k
BatchTopK	encoder.weight: shape (16,384, 2,304)		decoder.weight: shape (2,304, 16,384)		k
JumpReLU	W_enc: shape (2,304, 16,384)	b_enc: shape (16,384)	W_dec: shape (16,384, 2,304)	b_dec: shape (2,304)	threshold: shape (16,384)

## A.2 Training, Sparsity, and Compute Setup

Optimization and schedule. Adam with learning rate  $3\times10^{-4}$ ; LR warmup 1000 steps; sparsity warmup 5000 steps; LR decay starting at 80% of total steps. Precision: bfloat16. LM batch size = 4, context length = 2048, SAE batch size = 2048. Each run trains on  $\sim 5\times10^8$  tokens.

Sparsity controls. We sweep six target activity levels

$$L_0 \approx \{50, 80, 160, 320, 520, 820\}.$$

For TopK/BatchTopK we set k equal to the chosen  $L_0$  (aux-k coefficient 1/32; moving-threshold momentum 0.999; threshold tracking begins at step 1000). JumpReLU uses the same set via target\_10. For  $L_1$ -penalized families, we search the following penalty grids:

	Family	$L_1$ penalty values (used to span sparsity levels)				
550	Standard / Standard-New Gated SAE	0.012, 0.015, 0.020, 0.030, 0.040, 0.060 0.012, 0.018, 0.024, 0.040, 0.060, 0.080				

**Training details.** All training uses two NVIDIA RTX A800 GPUs. The table below reports the aggregated artifacts and training time (hours) for 30 SAEs per model (total 90), together with the runtime configuration. Times and sizes are approximate.

Model	#SAEs	Disk (GB)	Traing Time (H)	LM Batch	Context	SAE Batch	Peak Mem (GB)
Gemma- 54 2-2B	30	8.7	17	4	2048	2048	20
Gemma- 2-9B	30	13.2	60	4	2048	2048	70
Qwen- 2.5-3B	30	7.7	37	4	2048	2048	30

## **B** SAEBENCH Details, Results and Our Costs

## **B.1** Automated Interpretability Score Process

SAEBENCH [Karvonen et al., 2025] follow an LLM-as-judge pipeline to assign an *automated interpretability score* to each SAE latent. First, we collect layer activations by running the base LM with caching and encoding the residual stream through the SAE to obtain  $h \in \mathbb{R}^{N \times L \times F}$ . We define a token window of length 21 (buffer = 10) around any center (i,t) and, unless stated otherwise, mask BOS/PAD/EOS positions. For a latent  $\ell$ , we sample three window types: (i) **Top** (n=12 non-overlapping peaks of  $h[:,:,\ell]$ ), (ii) **Importance-Weighted** (n=7), sampled proportional to activation after removing values at least as large as the smallest Top peak), and (iii) **Random** (n=10), uniform over valid centers). Let  $v_{\max}$  be the maximum activation seen in any Top window position and set a global threshold  $\tau_{\text{act}} = 0.01 \, v_{\text{max}}$ .

We split the sampled windows into a **generation set** (10 Top + 5 IW) and a **scoring set** (2 Top + 2 IW + 10 Random, shuffled). In generation, tokens with activation  $> \tau_{\text{act}}$  are bracketed to highlight evidence; the judge LLM receives these 15 windows and returns a short English description of when the latent fires. In scoring, the judge sees the description and the 14 held-out windows without highlights and outputs a comma-separated list of indices it predicts as activations (or None).

Ground truth for a window  $\mathcal{W}$  is  $\mathbb{1}[\max_{u \in \mathcal{W}} h[u, \ell] > \tau_{\text{act}}]$ . The per-latent score is the **accuracy** over the M = 14 scoring windows, i.e.,

$$Score(\ell) = \frac{1}{M} \sum_{m=1}^{M} \mathbf{1} [\widehat{y}_m = y_m],$$

where  $\hat{y}_m \in \{0, 1\}$  is the judge prediction and  $y_m$  is the label defined above. For each SAE  $\theta$ , we evaluate 1,000 latents and report the mean over a random CONCEPT100 subset:

$$\mu(\theta) = \frac{1}{100} \sum_{\ell \in \text{CONCEPT100}} \text{Score}(\ell).$$

#### **B.2** Performance of SAEs on three models on SAEbench

Across the three backbones, the six SAEBENCH metrics (for information about these indicators, see SAEBENCH [Karvonen et al., 2025])jointly reveal how sparsity mechanisms balance interpretability, faithfulness, and causal structure. Automated Interpretability is strongest when encoders enforce compact latent usage (e.g., TopK/BatchTopK and ReLU at lower  $L_0$ ), and it gradually softens as capacity expands. The Absorption metric (considered via its complement in the plots) indicates that designs concentrating signal into a small set of latents are less prone to feature stealing, whereas higher effective capacity encourages redundancy and competition across latents. Meanwhile, Core/Loss-Recovered remains uniformly high, showing that even sparse codes closely preserve original model behavior; increasing  $L_0$  pushes faithfulness toward a ceiling without overturning the core trade-offs visible in the other metrics.

**Gemma-2-2B.** As shown in Fig. 5, Gemma-2-2B exhibits a balanced profile: interpretability stays robust for TopK/BatchTopK and ReLU at modest sparsity; absorption is contained when the code remains compact; and Core is near-saturated across the range. Improvements in SCR@20 are steady but measured, suggesting targeted debiasing with small k. Sparse Probing indicates that relatively few latents already carry much of the predictive signal, while RAVEL strengthens with moderate capacity, reflecting cleaner separation of attributes without undermining compactness.

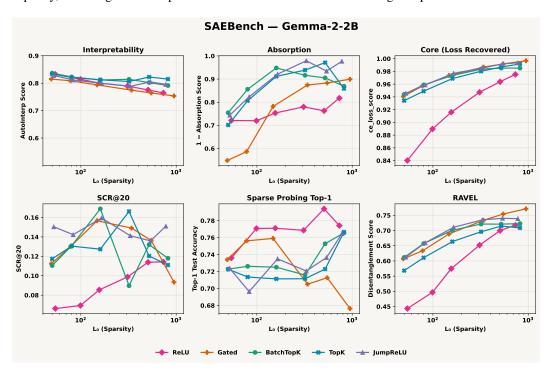


Figure 5: SAEbench results for **Gemma-2-2B**: interpretability remains strong at lower  $L_0$ , absorption stays low for compact codes, Core is near ceiling, and structure (SCR/RAVEL) improves with moderate capacity.

**Qwen-2.5-3B.** For Qwen-2.5-3B (Fig. 6), interpretability at low-to-moderate  $L_0$  is competitive—especially for TopK and JumpReLU—yet the model is more sensitive to absorption as capacity grows, implying greater latent competition and signal spread. Core remains excellent, so reconstructions are faithful; however, SCR gains can flatten at high  $L_0$  where residual spurious cues reappear. Sparse Probing is solid but a touch behind the strongest Gemma configurations, consistent with its flatter RAVEL patterns: causal structure is present but less crisply disentangled when attributes begin to diffuse across latents.

**Gemma-2-9B.** Gemma-2-9B (Fig. 7) pushes the upper envelope on structure: interpretability remains solid for compact encoders; absorption is low at moderate  $L_0$  that avoids unnecessary latent proliferation; and Core is near its ceiling. SCR@20 is the most decisive among the three, pointing to cleaner isolation of spurious factors with small, targeted ablations. Sparse Probing is strong and, together with higher RAVEL, indicates that only a handful of latents capture both predictive signal and causally specific attributes with minimal collateral interference.

## **B.3** SAEBench Runtime Cost

The computational requirements for running SAEBench evaluations were measured on two NVIDIA RTX A800 GPUs using **16K**-width SAEs trained on the Gemma-2-2B [Team et al., 2024], Qwen-2.5-3B [Yang et al., 2024] and Gemma-2-9B. Table 4 summarizes the *per-SAE* runtime for each evaluation type. Several evaluations include a one-time setup phase (e.g., precomputing activations or training probes) that can be reused across multiple SAEs; after this setup, each evaluation has its own runtime per SAE. We therefore report amortized per-SAE minutes.

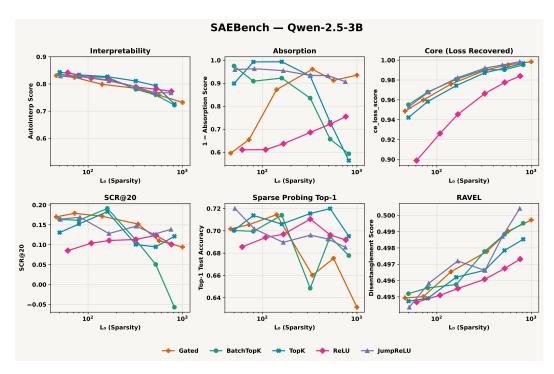


Figure 6: SAEbench results for **Qwen-2.5-3B**: strong interpretability at lower  $L_0$ , rising absorption with capacity, consistently high Core, and more fragile SCR/RAVEL at the highest capacities.

Table 4: **Approximate SAEBench runtime per SAE** (minutes). Values are per-SAE and represent amortized minutes after any one-time setup; each minute figure is an approximation and may vary with hardware and I/O.

Model	Core	Interpretability	Absorption	Sparse Probing	Ravel	SCR
Gemma-2-2B	4	8	12	2	18	10
Qwen-2.5-3B	7	9	15	8	17	16
Gemma-2-9B	11	12	17	30	40	28

## C AXBENCH Details and Our Costs

## C.1 Steering Score Process

$$\mathrm{HM}(C,I,F) = \begin{cases} 0, & \text{if } C{\le}0 \text{ or } I{\le}0 \text{ or } F{\le}0, \\ \frac{3}{\frac{1}{C} + \frac{1}{I} + \frac{1}{F}}, & \text{otherwise,} \end{cases} \in [0,2].$$

On the development split we pick  $\alpha^* = \arg\max_{\alpha \in \mathcal{A}} \frac{1}{5} \sum HM$ , then evaluate the held-out split with  $\alpha^*$  to obtain the *per-concept Steering Score* (mean HM over the 5 held-out prompts). SAE-level utility  $q(\theta)$  is the average of per-concept scores over the evaluation set for that SAE.

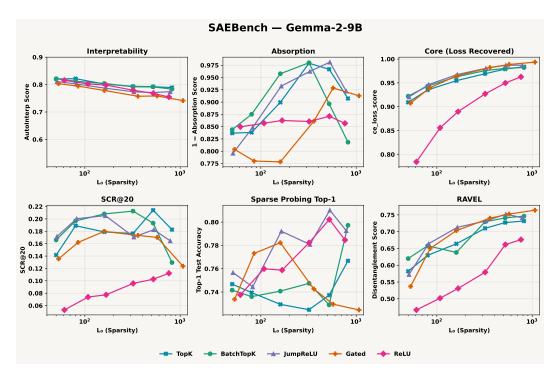


Figure 7: SAEbench results for **Gemma-2-9B**: robust interpretability with compact codes, low absorption at moderate  $L_0$ , near-ceiling Core, and the clearest gains in SCR/RAVEL among the three backbones.

## CONCEPT\_TEMPLATE

[System]

Please act as an impartial judge and evaluate whether the specified concept is presented in the sentence fragment provided below. Focus solely on whether the concept is clearly incorporated, without regard for grammar or logical coherence.

Begin your evaluation by providing a short explanation of whether the concept is clearly incorporated. Point out references to the concept in the fragment, noting any minimal or awkward incorporations. After providing your explanation, rate the concept's relevance on a scale from 0 to 2, where 0 indicates the concept is not present at all, 1 indicates the concept is somewhat present but minimally or awkwardly incorporated, and 2 indicates the concept is more fully and effectively incorporated, with stronger and more natural integration. Provide your rating using this exact format: "Rating: [[score]]".

[Concept Start] {concept} [Concept End]

[Sentence Fragment Start] {sentence} [Sentence Fragment End]

626

## INSTRUCT\_TEMPLATE

[System]

Please act as an impartial judge and evaluate whether the sentence fragment provided below is related to the instruction. Focus solely on the degree of relatedness in terms of topic, regardless of grammar, coherence, or informativeness.

Begin your evaluation by providing a brief explanation of whether the sentence is related to the instruction, and point out references related to the instruction. After providing your explanation, rate the instruction relevance on a scale from 0 to 2, where 0 indicates the sentence is unrelated to the instruction, 1 indicates it is somewhat related but only minimally or indirectly relevant in terms of topic, and 2 indicates it is more clearly and directly related to the instruction. Provide your rating using this exact format: "Rating: [[score]]".

[Instruction Start] {instruction} [Instruction End] [Sentence Fragment Start] {sentence} [Sentence Fragment End]

#### FLUENCY\_TEMPLATE

[System]

Please act as an impartial judge and evaluate the fluency of the sentence fragment provided below. Focus solely on fluency, disregarding its completeness, relevance, coherence with any broader context, or informativeness.

Begin your evaluation by briefly describing the fluency of the sentence, noting any unnatural phrasing, awkward transitions, grammatical errors, or repetitive structures that may hinder readability. After providing your explanation, rate the sentence's fluency on a scale from 0 to 2, where 0 indicates the sentence is not fluent and highly unnatural (e.g., incomprehensible or repetitive), 1 indicates it is somewhat fluent but contains noticeable errors or awkward phrasing, and 2 indicates the sentence is fluent and almost perfect. Provide your rating using this exact format: "Rating: [[score]]".

[Sentence Fragment Start] {sentence} [Sentence Fragment End]

628

629

## C.2 AxBench Steering Evaluation Cost

All steering-score evaluations were run on **two NVIDIA RTX A800 GPUs**. The LLM-as-judge backend was gpt-4o-mini. Evaluating one SAE on CONCEPT100 costs approximately \$5 in judge API fees; with **90** SAEs total ( $\approx 30$  per model), the per-model API cost is about \$150. Table 5 lists approximate per-SAE runtime and peak VRAM for each model.

Table 5: **AxBench steering evaluation cost per SAE.** Runtimes are per-SAE (hours) and approximate; VRAM is peak memory (GB). Judge fees assume gpt-4o-mini:  $\sim$  \$5 per SAE on CONCEPT100; Per-Model Cost assumes  $\sim$  30 SAEs/model ( $\approx$  \$150).

Model	Runtime / SAE (h)	Peak VRAM (GB)	Per-Model Cost (USD)
Gemma-2-2B	15	10	150
Qwen-2.5-3B	16	12	150
Gemma-2-9B	23	36	150

## **D** Implementation of $\Delta$ Token Confidence

For a fixed, neutral prefix s (we use "From my experience,", following the previous work[Arad et al., 2025]) we compare the next-token distribution of the base model with that of an intervened model in which a single SAE feature is amplified at layer L by a factor  $\alpha$ . The intervention is applied via the same SAE hook point used during training (on the residual stream of block L). We then compute the change in a confidence surrogate built from the top-k probabilities.

**Token confidence. [Fu et al., 2025]** For a distribution p over the vocabulary, let  $p_{(1)} \ge \cdots \ge p_{(k)}$  be the top-k probabilities.

$$C_k(p) = -\frac{1}{k} \sum_{i=1}^k \log p_{(i)}.$$

**Delta token confidence.** With  $p_{\text{base}}$  from a standard forward pass and  $p_{\text{int}}$  from a pass with the SAE feature intervention,

$$\Delta C_k(f; \alpha, L) = C_k(p_{\text{int}}) - C_k(p_{\text{base}}).$$

640 641

643

644

645

647

634

635

638

639

Each feature f is evaluated with two single-step forwards on the same prefix s: (i) a baseline pass; (ii) an intervened pass where we scale the code for f by  $\alpha$  before decoding it into the residual at layer L while keeping all other codes at zero. Hooks are removed immediately after the intervened pass to prevent accumulation across evaluations. In this work, we choose  $\alpha = 10$  and k = 1.

Feature selection from  $\Delta C_k$ . For each SAE we rank its features by  $\Delta C_k$  in two directions: UP (largest positive  $\Delta C_k$ ) and DOWN (most negative  $\Delta C_k$ ). We form selection sets using either (i) top-K by magnitude with  $K \in \{1, 2, 3, 4, 5\}$  per direction, or (ii) upper/lower-tail quantiles (e.g.,  $q \in \{0.99, 0.95, 0.90, 0.80\}$  mirrored for the lower tail). These sets are then carried into AXBENCH [Wu et al., 2025] to measure utility lift.

## 650 E Steering Results of SAE Architectures After Feature Selection

We quantify steering with the AXBENCH judge after selecting features using  $\Delta$  *Token Confidence* (Appendix D). Unless otherwise noted, lifts are reported as the percentage change of a given SAE's *steering score* relative to its own baseline (no selection). Results are organized at three levels: aggregate across SAEs per base model, per-SAE rankings, and distribution by architecture.

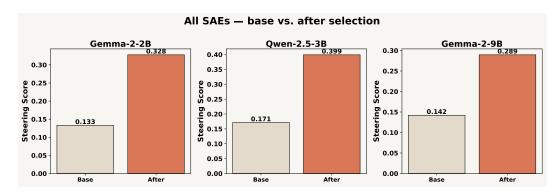


Figure 8: Overall steering score before and after feature selection. For each base model, the panel shows two bars: the average baseline steering score across its SAEs and the average after applying  $\Delta$  Token Confidence–based selection. Bars are annotated with the corresponding values; axes share the same scale across panels.

The aggregate view in Figure 8 summarizes how selection affects the mean *steering score* across all SAEs of a base model. Using  $\Delta$  *Token Confidence* for feature selection markedly improves the steering score across all three models in the figure. For Gemma-2-2B, the score rises from 0.133 to 0.328, which is a 146.6% relative improvement. Qwen-2.5-3B increases from 0.171 to 0.399, a 133.3% improvement, and achieves the highest post-selection score overall. Gemma-2-9B moves from 0.142 to 0.289, a 103.5% improvement.

**Conclusions:** (i) feature selection via  $\Delta$  *Token Confidence* consistently boosts steering for all models; (ii) relative gains are largest for the smallest model (Gemma-2-2B) and smallest for the largest model (Gemma-2-9B), suggesting diminishing relative returns with scale; and (iii) in absolute terms, Qwen-2.5-3B reaches the strongest final *steering score* after selection.

Figure 9 ranks SAEs within each model by their relative lift. Architecturally, no single SAE training approach dominates; however, the top-ranked lifts are frequently occupied by BatchTopK and Gated variants, with ReLU/JumpReLU also contributing strongly and TopK showing more mixed outcomes. Overall,  $\Delta$  Token Confidence yields consistent per-SAE gains, with variance decreasing and stability increasing as model size grows, while architectural diversity remains valuable for capturing the largest lifts.

Figure 10 groups lifts by architecture to visualize differences in central tendency and dispersion under the same selection and evaluation protocol. Read together with the per-SAE ranking, this distributional view helps disentangle architecture effects from model-specific variation and indicates which families tend to produce more stable or more variable outcomes after feature selection. BatchTopK and Gated generally occupy the highest central tendency with wide—but mostly positive—spread, especially on Gemma-2-2B and Qwen-2.5-3B. BatchTopK achieves the most stable and sizable steering gains. Variance is largest for the smallest model (Gemma-2-2B), indicating architecture-sensitive wins at small scale.

## F Dataset

Training corpus for SAEs. We train all SAEs on The Common Pile v0.1 [Kandpal et al., 2025], an openly licensed  $\sim$ 8 TB text collection built for LLM pretraining from  $\sim$ 30 sources spanning research papers, code, books, encyclopedias, educational materials, and speech transcripts. The corpus was curated as a principled alternative to unlicensed web text and has been validated by

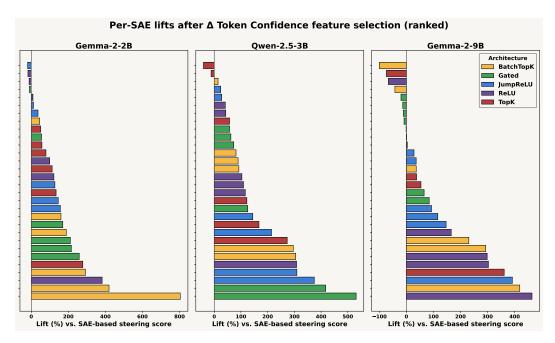


Figure 9: **Per-SAE percentage lift after**  $\Delta$  **Token Confidence selection.** Each panel corresponds to a base model. Horizontal bars report the percent lift of the SAE-level *steering score* relative to its own baseline, sorted from largest to smallest within the panel. Bar colors indicate the SAE training architecture (legend shared across panels).

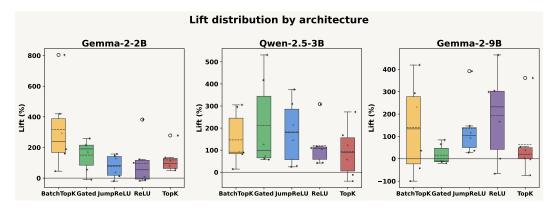


Figure 10: **Lift distributions by SAE architecture.** For each base model, box-and-whisker plots (with individual points overlaid) summarize the distribution of percentage lifts grouped by training architecture. Dashed horizontal lines denote the mean within each group, and whiskers follow the conventional interquartile rule.

training competitive 7B models on 1–2T tokens. We use it as the sole pretraining dataset for all SAE runs. More training details provided in Appendix A.2.

CONCEPT100 for steering utility. To evaluate steering, we construct CONCEPT100: a compact benchmark of 100 human-readable concept descriptions per evaluation set, produced automatically by our interpretability pipeline (Appendix B). Each entry is a pair (layer\_feature\_id, description) that summarizes a latent's activation pattern in plain language (e.g., mathematical symbols, scientific terms, pronouns, or domain phrases). These descriptions are supplied to the AXBENCH judge when computing steering score. The examples below illustrate the style and domain coverage.

## **Gemma-2-9B, BatchTopK,** $L_0 \approx 80$ . Ten examples:

- 20\_14429: concepts related to optical communication systems and their performance characteristics
- 20\_5795: specific technical terms and chemical compounds often related to scientific contexts
- 3. 20\_7908: terms related to gravitational lensing and its effects in cosmology
- 4. 20\_3042: pronouns and verbs indicating relationships or contributions in various contexts
- 20\_11897: scientific measurements and units related to energy, concentration, or biological data
- 6. 20\_12944: terms related to cell types and apoptosis mechanisms in scientific contexts
- 20\_8796: references to specific authors and statistical concepts in mathematical contexts
- 8. 20\_6430: the phrase "action" in mathematical and theoretical contexts
- 20\_2220: chemical elements and compounds, particularly including metals and metalrelated terms
- 10. 20\_585: various forms of the word "energy" and related concepts in scientific contexts

693

## **Qwen-2.5-3B, Gated,** $L_0 \approx 72$ **.** Ten examples:

- 1. 17\_15113: terms related to mathematical concepts and various scientific names or terms
- 2. 17\_11476: the phrase "as a function of" in contexts of measurement and analysis
- 3. 17\_162: mathematical symbols and concepts related to coordinates, magnitudes, and parameters in equations
- 17\_2552: dataset identifiers and technical terms common in research and academic documents
- 17\_16377: mathematical notation and technical terms commonly found in formal documents
- 6. 17\_3195: demographic, clinical, and biological characteristics in study populations and related comparisons
- 17\_9186: specific technical terms and concepts related to networking and programming
- 8. 17\_11487: mathematical notation and variables related to functions and equations
- 17\_14657: mathematical notations and structures involving angle brackets and properties of functions
- 10. 17\_1256: terms related to errors and error correction in coding theory and quantum operations

694

## **Gemma-2-2B, JumpReLU,** $L_0 \approx 81$ . Ten examples:

- 20\_11531: terms related to sports, programming, or specific keywords from various contexts
- 2. 20\_10460: terms related to fractional differential equations and numerical methods for solving them
- 3. 20\_4882: terms related to asymptotic theory, robustness, and statistical estimation methods
- 4. 20\_4425: first-person plural pronouns and expressions of intention or conjecture
- 5. 20\_372: technical terms related to measurement and structure in scientific contexts
- 6. 20\_9999: the word "from" and contexts implying deviation or distance from something

- 7. 20\_9703: the word "based" in various contexts of theoretical foundations and methodologies
- 8. 20\_15509: technical or numerical concepts in a variety of contexts
- 9. 20\_8218: phrases indicating conditions or assumptions that must be met in theoretical contexts
- 10. 20\_4614: time intervals and durations mentioned in the context of studies or observations

696 697 698

699

We currently maintain 90 SAEs (30 per base model). Beyond the CONCEPT100 sets evaluated in this paper, we have constructed the CONCEPT1000 and CONCEPT16K suites that scale the number of human-readable concepts up to 16K. We will extend training and evaluation to these larger suites in forthcoming releases to further substantiate the reliability and generality of this work.