

# EFFICIENT TRAINING OF SPARSE AUTOENCODERS FOR LARGE LANGUAGE MODELS VIA LAYER GROUPS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Sparse Autoencoders (SAEs) have recently been employed as an unsupervised approach for understanding the inner workings of Large Language Models (LLMs). They reconstruct the model’s activations with a sparse linear combination of interpretable features. However, training SAEs is computationally intensive, especially as models grow in size and complexity. To address this challenge, we propose a novel training strategy that reduces the number of trained SAEs from one per layer to one for a given group of contiguous layers. Our experimental results on Pythia 160M highlight a speedup of up to 3x without compromising the reconstruction quality and performance on downstream tasks. Therefore, layer clustering presents an efficient approach to train SAEs in modern LLMs.

## 1 INTRODUCTION

With the significant adoption of Large Language Model (LLM)s in world applications, understanding their inner workings has gained paramount importance. A key challenge in LLMs interpretability is the polysemanticity of neurons in models’ activations, lacking a clear and unique meaning (Olah et al., 2020). Recently, SAEs (Huben et al., 2024; Bricken et al., 2023) have shown great promise to tackle this problem by decomposing the model’s activations into a sparse combination of human-interpretable features.

The use of SAEs as an interpretability tool is motivated by two key reasons: the first is the substantial empirical evidence supporting the *Linear Representation Hypothesis (LRH)*, or that LLMs exhibit interpretable linear directions in their activation space (Mikolov et al., 2013; Nanda et al., 2023; Park et al., 2023); the second is the *Superposition Hypothesis (SH)* (Elhage et al., 2022), which supposes that, by leveraging sparsity, neural networks represent more features than they have neurons. Under this hypothesis, we can consider a trained neural network as a compressed simulation of a larger disentangled model, where every neuron corresponds to a single feature. To overcome superposition, SAEs leverage the LRH to decompose model activations into a sparse linear combination of interpretable features.

However, training SAEs is computationally expensive and will become even more costly as the model size and parameter counts grow. Indeed, one Sparse Autoencoder (SAE) is typically learned for a given component at every layer in a LLM. Moreover, the number of features usually equals the model activation dimension multiplied by a positive integer, called the expansion factor. For example, a single SAE trained on the Llama-3.1 8B model activations with an expansion factor of 32 has roughly  $4096^2 \times 32 \times 2 \approx 1.073$  billion parameters, resulting in a total of more than 32 billions parameters when training one for each of the 32 layers. Additionally, every SAE is trained on a fixed number of tokens  $T$ , typically in the billions, leading to a cumulative training requirement of  $L \times T$  tokens, where  $L$  is the total number of layers—for the Llama-3.1 8B model, this amounts to 32 billion tokens. This high computational demand not only increases training time but also requires substantial hardware resources and energy consumption, making the approach increasingly impractical as models scale.

In this work, we reduce the computational overhead of training a separate SAE for each layer of a target LLM by learning a single SAE for groups of related and contiguous layers, thus reducing

---

\*Equal contribution

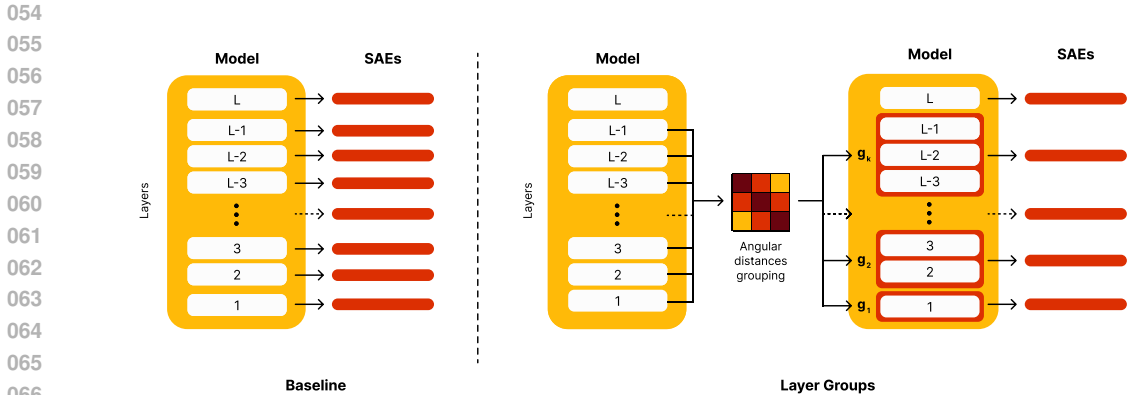


Figure 1: The illustration of our method. While standard training of SAEs requires training one per layer, our method first clusters layers by angular similarity and then trains a single SAE for each group.

the number of trained SAEs from  $L$  to  $k$ , where  $k$  is the number of groups. Furthermore, to ensure computational efficiency remains consistent, we fix the total number of training tokens  $T$  per group and allocate  $T/l$  tokens to each layer within the group, where  $l$  is the number of layers in that group. Thus, our approach reduces the total number of training tokens from  $LT$  to  $kT$ . This approach, which we term **Group-SAE**, is inspired by the observation that neural network layers often group together in learning task-specific representations (Szegedy et al., 2014; Zeiler & Fergus, 2014; Jawahar et al., 2019): shallow layers typically capture low-level features, while deeper layers learn high-level abstractions. Additionally, adjacent layers in LLMs tend to encode redundant information, as evidenced by the similarity in the angular distance of their outputs (Gromov et al., 2024). To quantify the relatedness between layers, we measure the angular distance of the residual stream after the MLP’s contribution across neighboring layers.

Denoting with  $L$  the number of layers of the target model and  $k$  the number of groups we want to cluster the layers into, our approach obtains at least a  $\frac{(L-1)}{k}$  times speedup<sup>1</sup> without sacrificing neither reconstruction quality nor downstream performance.

Our **contributions** can be summarized as follows:

- We demonstrate that a single Group-SAE can effectively reconstruct the activations of an entire cluster as a sparse linear combination of interpretable features. This approach reduces the number of SAEs required to be trained for a given model by a factor of  $k$ .
- We demonstrate the practical effectiveness of our method by training SAEs on the Pythia-160M model for different values of  $k$ , and we obtain a  $\frac{L-1}{k}$  times speedup at the cost of a minor deterioration in performance.
- We extensively analyze our method on reconstruction performance, downstream tasks, and human interpretability for different values of  $k$ .

## 2 RELATED WORK

### 2.1 THE LINEAR REPRESENTATION AND SUPERPOSITION HYPOTHESES

Supported by substantial evidence, from the seminal Mikolov et al. (2013) vector arithmetic to the more recent work of Nanda et al. (2023) and Park et al. (2023) on LLMs, the Linear Representation Hypothesis (LRH) supposes that neural networks have interpretable linear directions in their acti-

<sup>1</sup>We do not consider the last transformer layer, as different from every other layer w.r.t. the angular distance defined in Section 3.1.

variations space. However, neuron polysemanticity remains an essential challenge in neural network interpretability (Olah et al., 2020).

Recently, Bricken et al. (2023) explored this issue by relating the Superposition Hypothesis (SH) to the decomposition ideally found by a SAE. According to the SH, neural networks utilize  $n$ -dimensional activations to encode  $m \gg n$  features by leveraging their sparsity and relative importance. As a result, we can write the activations  $\mathbf{x}^j$  in a model as

$$\mathbf{x}^j \approx \mathbf{b} + \sum_i^m \mathbf{f}_i(\mathbf{x}^j) \mathbf{d}_i \quad (1)$$

where  $\mathbf{x}^j \in \mathbb{R}^n$  is the activation vector for an example  $j$ ,  $\mathbf{f} \in \mathbb{R}^m$  is a sparse feature vector,  $\mathbf{f}_i(\mathbf{x}^j)$  is the activation of the  $i$ -th feature,  $\mathbf{d}_i \in \mathbb{R}^n$  is a unit vector in the activation space and  $\mathbf{b} \in \mathbb{R}^n$  is a bias term.

## 2.2 SPARSE AUTOENCODERS

Sparse Autoencoders have gained popularity in LLM interpretability due to their ability to counteract superposition and decompose neuron activations into interpretable features (Bricken et al., 2023; Huben et al., 2024). Given an input activation  $\mathbf{x} \in \mathbb{R}^{d_{\text{model}}}$ , a SAE reconstructs it as a sparse linear combination of  $d_{\text{sae}} \gg d_{\text{model}}$  features, denoted as  $\mathbf{v}_i \in \mathbb{R}^{d_{\text{model}}}$ , where  $d_{\text{sae}}$  is set as  $d_{\text{sae}} = c \cdot d_{\text{model}}$ , where  $c \in \{2^n | n \in \mathbb{N}_+\}$ . The reconstruction follows the form:

$$(\hat{\mathbf{x}} \circ \mathbf{f})(\mathbf{x}) = \mathbf{W}_d \mathbf{f}(\mathbf{x}) + \mathbf{b}_d \quad (2)$$

Here, the columns of  $\mathbf{W}_d \in \mathbb{R}^{d_{\text{sae}} \times d_{\text{model}}}$  represent the features  $\mathbf{v}_i$ ,  $\mathbf{b}_d \in \mathbb{R}^{d_{\text{model}}}$  is the decoder’s bias term, and  $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^{d_{\text{sae}}}$  represents the sparse features activations. The feature activations are computed as

$$\mathbf{f}(\mathbf{x}) = \sigma(\mathbf{W}_e(\mathbf{x} - \mathbf{b}_d) + \mathbf{b}_e) \quad (3)$$

where  $\mathbf{W}_e \in \mathbb{R}^{d_{\text{sae}} \times d_{\text{model}}}$  and  $\mathbf{b}_e \in \mathbb{R}^{d_{\text{sae}}}$  are the encoder’s matrix and bias term respectively, and  $\sigma$  is an activation function, typically  $\text{ReLU}(\mathbf{x}) = \max(0, \mathbf{x})^2$ . The training of a SAE involves minimizing the following loss function:

$$\mathcal{L}_{\text{sae}} = \|\mathbf{x} - \hat{\mathbf{x}}(\mathbf{f}(\mathbf{x}))\|_2^2 + \lambda \text{S}(\mathbf{f}(\mathbf{x})) \quad (4)$$

where the first term in Equation 4 represents the reconstruction error, while the second term  $\text{S}$  is the regularization on the activations  $\mathbf{f}(\mathbf{x})$  to encourage sparsity, which is typically equal to the  $\ell_1$  norm of the features  $\mathbf{f}(\mathbf{x})$ , i.e.  $\text{S}(\mathbf{f}(\mathbf{x})) = \|\mathbf{f}(\mathbf{x})\|_1$  Bricken et al. (2023). Other definitions of sparsity are admissible: Rajamanoharan et al. (2024) set  $\text{S}$  equal to the  $\ell_0$  norm of the features, i.e.  $\text{S}(\mathbf{f}(\mathbf{x})) = \|\mathbf{f}(\mathbf{x})\|_0 = \sum_{i=0}^{d_{\text{sae}}} \mathbb{I}[\mathbf{f}_i(\mathbf{x}) \neq 0]$ , which is the one adopted in this work along with the JumpReLU activation function, while Gao et al. (2024) imposes sparsity by selecting the Top-K activating features from  $\mathbf{f}(\mathbf{x})$ .

## 2.3 SAES EVALUATION

SAE evaluation in the context of LLMs presents a significant challenge. While standard unsupervised metrics such as  $L_2$  (reconstruction) loss and  $L_0$  sparsity are widely adopted to measure SAE performance (Gao et al., 2024; Lieberum et al., 2024), they fall short of assessing two key aspects: causal importance and interpretability.

Recent approaches, including auto-interpretability (Bricken et al., 2023; Huben et al., 2024; Bills et al., 2023) and ground-truth comparisons (Sharkey et al., 2023), aim to provide a more holistic evaluation. These methods focus on the causal relevance of features (Marks et al., 2024) and evaluate SAEs in downstream tasks. Makelov et al. (2024), for instance, proposed a framework for the

<sup>2</sup>Other activation functions are admissible, e.g. Top-K (Gao et al., 2024) or JumpReLU (Rajamanoharan et al., 2024), with JumpReLU being the one employed in this work.

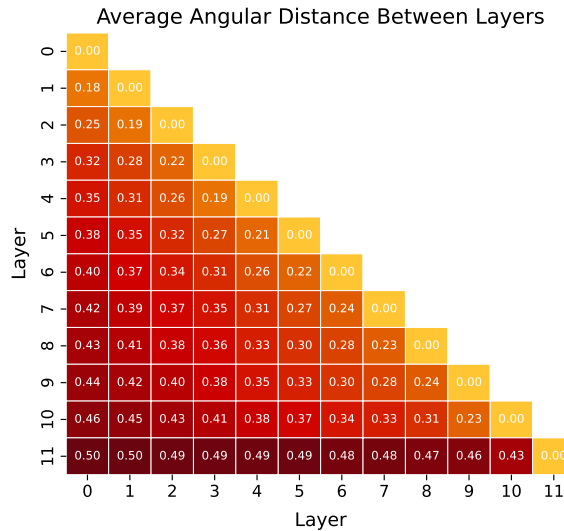


Figure 2: Average angular distance between all layers of the Pythia-160M model, as defined in Equation 5. The angular distances are computed over 5M tokens from the training dataset. The angular distances are bounded in  $[0, 1]$ , where an angular distance equal to 0 means equal activations, 0.5 means activations are perpendicular and an angular distance of 1 means that the activations point in opposite directions.

Indirect Object Identification (IOI) task, emphasizing three aspects: the sufficiency and necessity of reconstructions, sparse feature steering (Templeton et al., 2024), and the interpretability of features in causal terms.

Karvonen et al. (2024) further contributed by developing specialized metrics for board game language models. Using structured games like chess and Othello, they introduced supervised metrics, such as board reconstruction accuracy and coverage of predefined state properties, offering a more direct assessment of SAEs’ ability to capture semantically meaningful and causally relevant features.

## 2.4 IMPROVING SAEs TRAINING

As SAEs gain popularity for LLMs interpretability and are increasingly applied to state-of-the-art models (Lieberum et al., 2024), the need for more efficient training techniques has become evident. To address this, (Gao et al., 2024) explored the scaling laws of Autoencoders to identify the optimal combination of size and sparsity.

Recent work also explored using transfer learning to improve SAE training. For example, Kissane et al. (2024) and Lieberum et al. (2024) demonstrated the transferability of SAE weights between base and instruction-tuned versions of Gemma-1 (Team et al., 2024a) and Gemma-2 (Team et al., 2024b), respectively. On the other hand, Ghilardi et al. (2024) shows that transfer also occurs among layers of a single model, both in forward and backward directions.

## 3 EXPERIMENTAL SETUP

### 3.1 LAYER GROUPS

For a model with  $L$  layers, the number of possible combinations of  $k$  groups of adjacent layers that can be tested is given by  $\binom{L-1}{k-1}$ . With this number growing with model depth, we employed an agglomerative grouping strategy based on the angular distances between layers to reduce it drastically. In particular, we compute the mean angular distances, as specified in Gromov et al. (2024), over 5M tokens from our training set:

$$d_{\text{angular}}(\mathbf{x}_{\text{post}}^p, \mathbf{x}_{\text{post}}^q) = \frac{1}{\pi} \arccos\left(\frac{\mathbf{x}_{\text{post}}^p \cdot \mathbf{x}_{\text{post}}^q}{\|\mathbf{x}_{\text{post}}^p\|_2 \|\mathbf{x}_{\text{post}}^q\|_2}\right) \quad (5)$$

for every  $p, q \in \{1, \dots, L\}$ , where  $\mathbf{x}_{\text{post}}^l$  are the  $l$ -th residual stream activations after the MLP’s contribution<sup>3</sup>. From Figure8, it can be noted how the last layer is different from every other layer in terms of angular distance. For this reason, we have decided to exclude it from the grouping procedure.

We adopted a bottom-up hierarchical clustering strategy with a complete linkage (Nielsen, 2016) that aggregates layers based on their angular distances. Specifically, at every step of the process, the two groups with minimal group-distance<sup>4</sup> are merged; we repeat the process until the predefined number of groups  $k$  is reached. This approach prevents the formation of groups as long chains of layers and ensures that the maximum distance within each group remains minimal. In this work, we create groups considering all layers except the last, and we choose  $k$  varying from 1 (a single cluster with all layers) to 5. Layer groups can be found in Appendix B. To show how our method scales to larger models, we also compute distances and groups of Gemma-2 2b and 9b Team et al. (2024b). Detailed results are provided in Appendix B.

### 3.2 SAES TRAINING AND EVALUATION

We denote  $\text{SAE}_i$  as the baseline SAE trained to reconstruct the activations of layer  $i$ . For every  $1 \leq j \leq k \leq 5$  with  $j, k \in \mathbb{N}$ , we define  $\text{SAE}_{j_k}$  as the SAE trained to reconstruct the activations for all layers in the  $j$ -th group of a partition consisting of  $k$  groups. Furthermore, let  $[j_k]$  represent the set of layers belonging to the  $j$ -th group within this partition. Each  $\text{SAE}_{j_k}$  is specifically trained to reconstruct the activations of each layer  $s \in [j_k]$  individually, using the formulations described in equations 3 and 2.

To ensure both computational efficiency and fair comparison with baselines, we allocate a fixed total of 1 billion tokens for training both every  $\text{SAE}_i$  and  $\text{SAE}_{j_k}$ . In particular, for  $\text{SAE}_{j_k}$  these tokens are evenly distributed across the layers in  $[j_k]$ , assigning  $1\text{B}/|[j_k]|$  tokens to each layer, with a randomly chosen layer’s activation per token. This approach ensures balanced training across layers within the group while keeping the total token budget consistent with the baseline configurations, thus reducing the total number of training tokens from  $LT$  to  $kT$ .

For evaluation, we train a  $\text{SAE}_{j_k}$  for every combination of  $k$  and  $j$  and compare its performance with the corresponding baseline  $\text{SAE}_s$ , where  $s \in [j_k]$ . Section 4 presents the results of these comparisons based on standard reconstruction and sparsity metrics. Furthermore, in Section 5, we show how our approach performs on popular downstream tasks (e.g., (Marks et al., 2024; Hanna et al., 2023; Wang et al., 2023)), while Section 6 reports the human interpretability scores.

### 3.3 DATASET AND HYPERPARAMETERS

We train SAEs on the residual stream of the Pythia-160M model (Biderman et al., 2023) for 1B tokens, focusing on the residual stream activations after the MLP contribution. The chosen dataset is a 2B pre-tokenized version<sup>5</sup> of the Pile dataset (Gao et al., 2020) with a context size of 1024. We set the expansion factor  $c = 8$ ,  $\lambda = 3e-4$  in Equation 4, learning rate equal to  $7e-4$ , and a batch size of 4096 samples. Following Bricken et al. (2023), we constrain the decoder columns to have unit norm and do not tie the encoder and decoder weights. Additionally, we have normalized the activations so that they have mean squared  $\ell_2$  norm of one during SAE training, as specified in Rajamanoharan et al. (2024), estimating the scaling factor over 2M tokens of our train set.

We use the JumpReLU activation function as specified in Rajamanoharan et al. (2024), and defined as  $\text{JumpReLU}_\theta(z) = z \cdot \text{H}(z - \theta)$ , with  $\theta$  being a threshold learned during training and H is the

<sup>3</sup>In a Transformer model (Vaswani et al., 2017), the residual stream activations are computed as follows: given the activations  $\mathbf{x}_{\text{pre}}^{l-1}$  from layer  $l-1$ , the pre-layer activations for layer  $l$ ,  $\mathbf{x}_{\text{pre}}^l = \mathbf{x}_{\text{post}}^{l-1}$ , are calculated as  $\mathbf{x}_{\text{post}}^{l-1} = \mathbf{x}_{\text{mid}}^{l-1} + \text{MLP}(\mathbf{x}_{\text{mid}}^{l-1})$ , where  $\mathbf{x}_{\text{mid}}^{l-1} = \mathbf{x}_{\text{pre}}^{l-1} + \text{Attn}(\mathbf{x}_{\text{pre}}^{l-1})$ . All SAEs are trained specifically on the residual stream after the MLP contribution, i.e.,  $\mathbf{x}_{\text{post}}^l$ , for every layer  $l$ .

<sup>4</sup>The complete linkage clustering strategy defines the group-distance between two groups  $X$  and  $Y$  as  $D(X, Y) = \max_{\mathbf{x} \in X, \mathbf{y} \in Y} d_{\text{angular}}(\mathbf{x}, \mathbf{y})$

<sup>5</sup><https://huggingface.co/datasets/NeelNanda/pile-small-tokenized-2b>

Heaviside function. We fixed the hyperparameters for all the experiments conducted in this work. All hyperparameters can be found in Table 2.

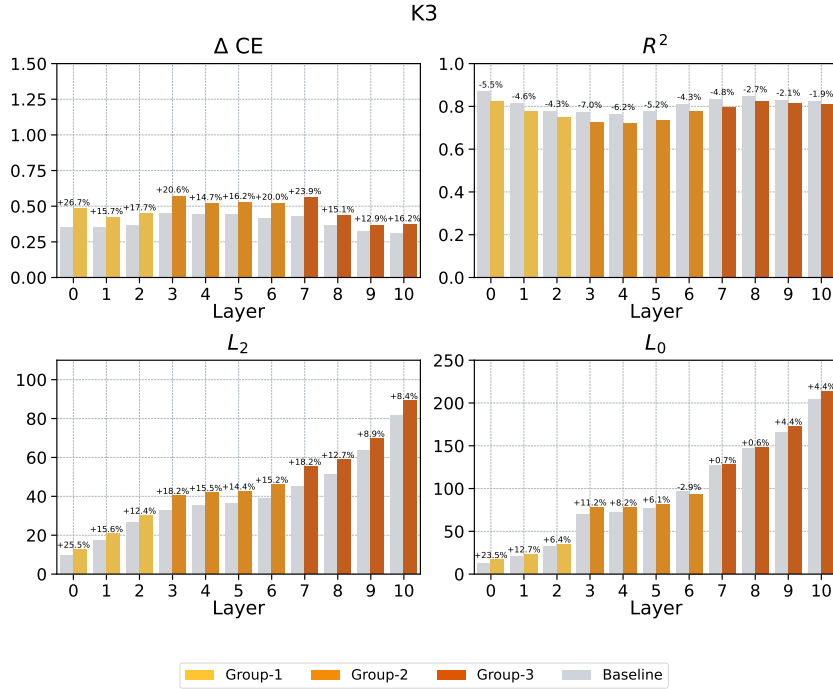


Figure 3:  $\Delta CE$ ,  $R^2$ ,  $L_2$  and  $L_0$  for  $k = 3$  number of groups. Note that, for  $\Delta CE$  and  $L_2$ , the lower is the score, the better it is, while the contrary is true for  $R^2$ . For  $L_0$  there are no indications on a preferred value.

#### 4 RECONSTRUCTION EVALUATION

To assess the quality of the reconstruction of SAEs trained with our grouping strategy, we report three standard reconstruction metrics and one sparsity metric. In particular, the Delta Cross-Entropy, defined as  $CE(\hat{x} \circ f) - CE(M)$ , measures the difference in Cross-Entropy loss (CE) between the output with SAE reconstructed activations ( $\hat{x} \circ f$ ) and the model’s output ( $M$ ).

The  $L_2$  loss is the first term of Equation 4 and measures the reconstruction error made by the SAE. The  $R^2$  score, defined as  $1 - \frac{\|x - \hat{x}\|_2^2}{\|x - \mathbb{E}_{x \sim D}[x]\|_2^2}$ , measures the fraction of explained variance of the input recovered by the SAE.

Finally, the  $L_0$  sparsity loss, defined as  $\sum_{j=1}^{d_{\text{sac}}} \mathbb{I}[f_j \neq 0]$ , represents the number of non-zero SAE features used to compute the reconstruction. For each metric, we compute the average over 1M examples from the test dataset.

Figure 3 shows the reconstruction and sparsity metrics for each layer of the grouping with  $k = 3$ . It can be noted that training a single SAEs on the activations from multiple close layers doesn’t dramatically affect the reconstruction, even when all layers are clustered in a small number of groups. Metrics for all other values of  $k$  can be found in Appendix C.

These results demonstrate that a single SAE can effectively reconstruct activations across multiple layers. Furthermore, the comparable performance between  $SAE_{j_k}$  and the individual layer-specific  $SAE_i$  indicates that post-residual stream activations from adjacent layers share a common set of underlying features. This hypothesis is further supported by directly comparing the directions learned by  $SAE_i$  and  $SAE_{j_k}$  using the Mean Maximum Cosine Similarity (MMCS) score, as shown in Appendix D.

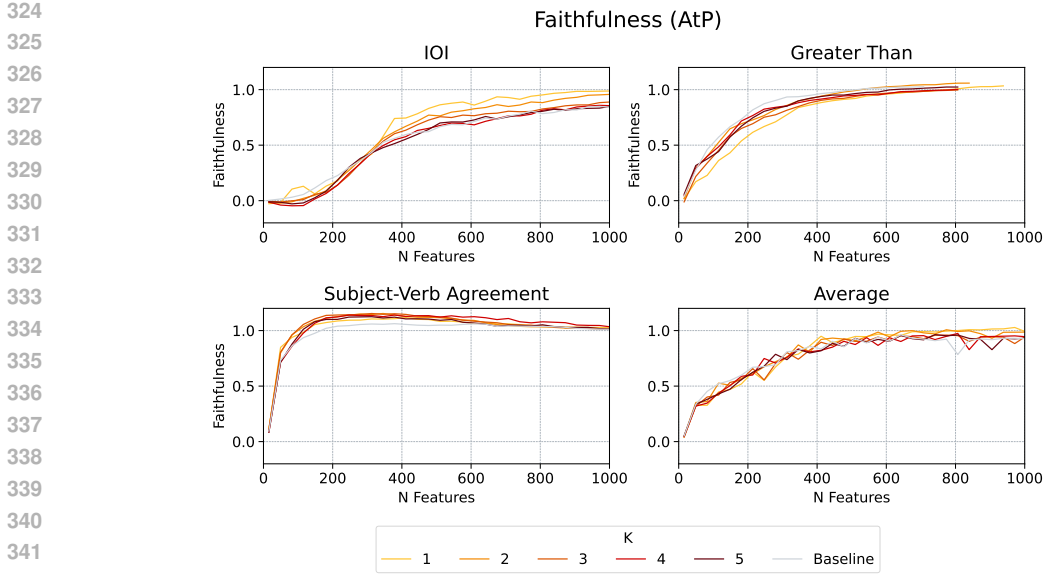


Figure 4: Average faithfulness (Marks et al., 2024) for every downstream task (IOI, Greater Than, Subject-Verb Agreement) with IE computed with AtP (Equation 8). The “Baseline” average is computed considering the performance obtained by  $SAE_i, \forall i = 0, \dots, 10$ . The “Average” plot depicts the average over the three downstream tasks.

## 5 DOWNSTREAM TASKS EVALUATION

While achieving good reconstruction metrics is crucial for a trained SAE, it is insufficient for a comprehensive evaluation of its performance. For instance, unsupervised metrics alone cannot determine whether the identified features capture causally influential directions for the model. To address this, following Marks et al. (2024), we applied SAEs to three well-known tasks: Indirect Object Identification (IOI), Greater Than, and Subject-Verb Agreement.

Each task can be represented as a set of counterfactual prompts paired with their respective answers, formally denoted as  $\mathcal{T} : \{x_{\text{clean}}, x_{\text{corrupted}}, a_{\text{clean}}, a_{\text{corrupted}}\}$ . Counterfactual prompts are similar to clean ones but contain slight modifications that result in a different predicted answer.

Ideally, a robust SAE should be able to recover the model’s performance on a task when reconstructing its activations. Furthermore, we expect the SAE to rely on a small subset of task-relevant features to complete the task. To assess this, we filtered the features to include only the most important ones, where importance is defined as the *indirect effect* (IE) (Pearl, 2022) of the feature on task performance, measured by a real-valued metric  $m : \mathbb{R}^{d_{\text{vocab}}} \rightarrow \mathbb{R}$ . Specifically,

$$IE(m; \mathbf{f}_i; x_{\text{clean}}, x_{\text{corrupted}}) = m(x_{\text{clean}} | \text{do}(\mathbf{f}_i = \mathbf{f}_{i;\text{corrupted}})) - m(x_{\text{clean}}) \tag{6}$$

In this equation,  $\mathbf{f}_{i;\text{corrupted}}$  represents the value of the  $i$ -th feature  $\mathbf{f}_i$  during the computation of  $m(x_{\text{corrupted}})$ , and  $m(x_{\text{clean}} | \text{do}(\mathbf{f}_i = \mathbf{f}_{i;\text{corrupted}}))$  refers to the value of  $m$  for  $x_{\text{clean}}$  under an intervention where the activation of feature  $\mathbf{f}_i$  is set to  $\mathbf{f}_{i;\text{corrupted}}$ . Moreover,  $m$  is defined as the difference in logits between the clean and the corrupted answer.

Calculating these effects is computationally expensive, as it requires a forward pass for each feature. To mitigate this, we employed two approximate methods: Attribution Patching (AtP)(Nanda, 2023; Syed et al., 2024) and Integrated Gradients (IG)(Sundararajan et al., 2017). Appendix E provides a formal definition of both methods.

Following Marks et al. (2024), we used faithfulness and completeness metrics to evaluate the performance of the SAEs on the tasks. These metrics are defined as  $\frac{m(C) - m(\emptyset)}{m(M) - m(\emptyset)}$ , where  $m(M)$  and  $m(\emptyset)$  represent the metric average over  $\mathcal{T}$ , achieved by the model alone and with the mean-ablated

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

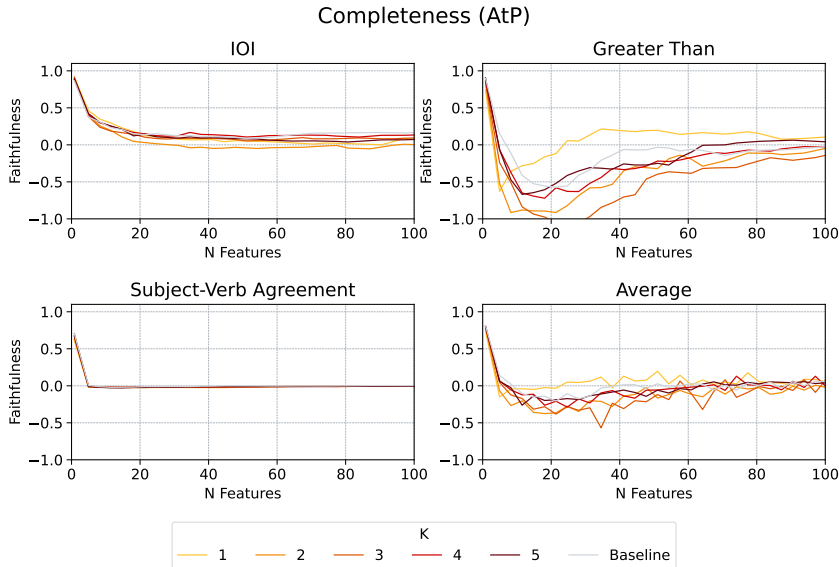


Figure 5: Average completeness for every downstream task (IOI, Greater Than, Subject-Verb Agreement) with IE computed with AtP (Equation 8). The “Baseline” average is computed considering the performance obtained by  $SAE_i, \forall i = 0, \dots, 10$ . The “Average” plot depicts the average performance over the three downstream tasks.

SAE reconstructions, respectively.  $m(C)$  is computed based on the task and either the *faithfulness* or *completeness* criteria: for faithfulness, it is the metric average when using only the important SAE features, while mean-ablating the others; on the contrary for completeness, it is calculated by mean-ablating the important features while keeping the others active.

Given a predefined number  $N$  of features to compute a circuit, we select the top- $N$  features based on their indirect effect at each layer. Then, to calculate the Faithfulness and Completeness scores, we respectively retain or ablate these selected features while mean-ablating or retaining all the others. As in Marks et al. (2024), the computation of both scores incorporates the SAE error  $\epsilon_x$ <sup>6</sup> into the reconstructions, as excluding it compromises the model’s performance on the task. Furthermore, when using Group-SAEs, the selected features can vary across layers within a group, as each layer may require a distinct set of features to achieve optimal reconstruction, thus  $SAE_{j_k}$  is used to independently and separately reconstruct the activations of every layer  $s \in [j_k]$  with the *indirect effect* that is estimated for every layer separately.

These metrics allow us to evaluate two critical aspects of SAE quality: whether the SAE learned a set of features that is both **sufficient** and **necessary** to perform the task. Figures 4 and 5 display the faithfulness and completeness scores for all  $k$  groups.

All  $SAE_{j_k}$  perform comparably to, or slightly better than, the baseline  $SAE_i$  models on the faithfulness score (Figure 4) across all three downstream tasks evaluated, demonstrating the sufficiency of learned features. Remarkably, this performance is achieved using only the top 5% of the most important active features, which recover 75% of the baseline performance, on average. Even more remarkable are the performances of the single-group  $SAE_{j_1}$  ( $k = 1$ ), which closely follow the trend of both the baseline SAEs and the  $SAE_{j_5}$  ( $k = 5$ ). Moreover, we did not observe any substantial differences in performance for any of the tested values of  $k$ . The necessity of the features learned by both the baseline  $SAE_i$  and  $SAE_{j_k}$  is confirmed by the completeness scores depicted in Figure 5, with a severe drop in performance even with only the top 10 active features mean-ablated.

The results of  $SAE_{j_k}$  on downstream tasks demonstrate that their learned features are both sufficient and necessary. Moreover, these findings confirm those in Section 4, i.e., that a single SAE can

<sup>6</sup>The SAE error  $\epsilon$  is defined as  $\epsilon_x = \mathbf{x} - \hat{\mathbf{x}}(\mathbf{x})$ .



effectively reconstruct activations across multiple contiguous layers and learn a set of shared, general features that span adjacent layers.

## 6 HUMAN INTERPRETABILITY

In addition to achieving excellent reconstruction and downstream performance, SAEs must learn interpretable features. Following Ghilardi et al. (2024), we engaged human annotators to identify interpretable patterns in the feature annotations. Specifically, they attempted to provide clear definitions for each feature by examining its top and bottom logits attribution scores, as well as the top activating tokens. In total, we randomly sample 96 features for each SAE and store their activations on 1M tokens from the training dataset. To ensure learned features are not layer-specific, for Group-SAEs we keep the sampled features constant across the layers within the group. To evaluate the quality of these features, we defined the *Human Interpretability Score* as the ratio of features considered interpretable by the human annotators.

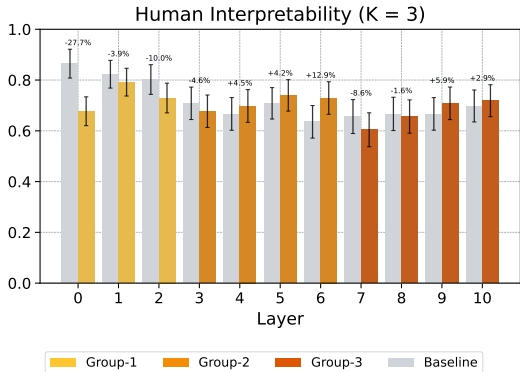


Figure 6: Human Interpretability Scores for  $k = 3$ . The differences in the interpretability scores of features learned by the  $SAE_{j_k}$  and the baseline  $SAE_i$  are not statistically significant different for all the layers except for the first. Error bars shows one standard deviations of the scores differences, modeled as a Binomial distribution (Wasserman, 2010).

Figure 6 presents the Human Interpretability Scores for all layers of  $k = 3$ . Scores for all other values of  $k$  can be found in Appendix F. According to human annotators, the interpretability of the features learned by  $SAE_{j_k}$  is comparable to that of the baseline  $SAE_i$ . Moreover, we found that when many layers are grouped together, e.g., when  $k = 1$ ,  $SAE_{j_k}$  features are more polysemantic overall, probably due to increased interference. Nevertheless, some of them remain perfectly interpretable across all model layers, capturing critical directions in model computations.

## 7 CONCLUSION

This work introduces a novel approach to efficiently train Sparse Autoencoders (SAEs) for Large Language Models (LLMs) by clustering layers based on their angular distance and training a single SAE for each group. Through this method, we achieved up to a 3x speedup in training without compromising reconstruction quality or performance on downstream tasks. The results demonstrate that activations from adjacent layers in LLMs share common features, enabling effective reconstruction with fewer SAEs.

Our findings also show that the SAEs trained on grouped layers perform comparably to layer-specific SAEs in terms of reconstruction metrics, faithfulness, and completeness on various downstream tasks. Furthermore, human evaluations confirmed the interpretability of the features learned by our SAEs, underscoring their utility in disentangling neural activations.

The methodology proposed in this paper opens avenues for more scalable interpretability tools, facilitating deeper analysis of LLMs as they grow in size. Future work will focus on further optimizing the number of layer groups and scaling the approach to even larger models.

## 8 LIMITATIONS AND FUTURE WORKS

One limitation of our approach is the absence of a precise method for selecting the optimal number of layer groups ( $k$ ). This choice is due to the lack of a clear *elbow rule* for identifying the correct number of groups.

Based on our results on Pythia-160m, Appendix B provides the Maximum Average Angular Distance (MAAD) within groups for Pythia-160m and larger models (Gemma-2 2b and 9b). This heuristic can guide the selection of the number of groups  $k$  for clustering layers in larger models, relying on the assumption that angular distance between activations from different layers is an important indicator of how the Group SAE will perform on them. While further investigation is needed to refine group configurations and assess scalability, these heuristics offer a practical reference for training Group-SAEs in larger models.

Additionally, we tested our method primarily on the Pythia-160M model, a relatively small LLM. While our findings demonstrate significant improvements in efficiency without sacrificing performance, the scalability of our approach to much larger models remains an open question. Future work could explore how the grouping strategy and training techniques generalize to models with billions of parameters, where the computational benefits would be even more pronounced.

Another important direction for future research involves understanding how Sparse Autoencoders (SAEs) handle the superposition hypothesis when encoding information from multiple layers. While our method effectively grouped layers and maintained high performance, how SAEs manage the potential overlap in feature representation across layers remains unclear. Investigating this aspect could lead to a more clear understanding of the trade-offs between sparsity and feature disentanglement in SAEs, and inform strategies for improving interpretability without compromising task performance.

In summary, while our work represents an efficient step forward in training SAEs for interpretability, extending this approach to larger models and exploring the handling of superposition will provide valuable insights for both practical applications and the theoretical understanding of sparse neural representations.

## 9 REPRODUCIBILITY STATEMENT

To support the replication of our empirical findings on training SAEs via layer groups and to enable further research on understanding their inner works, we plan to release all the code and SAEs used in this study upon acceptance.

## REFERENCES

- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. Language models can explain neurons in language models, 2023. URL <https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html>. Accessed: 2024-08-18.
- 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://transformer-circuits.pub/2023/monosemantic-features/index.html>.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish,

- 540 Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of super-  
541 position. *Transformer Circuits Thread*, 2022. URL [https://transformer-circuits.  
542 pub/2022/toy\\_model/index.html](https://transformer-circuits.pub/2022/toy_model/index.html).
- 543
- 544 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason  
545 Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The Pile:  
546 An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 547 Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever,  
548 Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders, 2024. URL [https:  
549 //arxiv.org/abs/2406.04093](https://arxiv.org/abs/2406.04093).
- 550 Davide Ghilardi, Federico Belotti, Marco Molinari, and Jaehyuk Lim. Accelerating sparse autoen-  
551 coder training via layer-wise transfer learning in large language models. In *The 7th BlackboxNLP  
552 Workshop*, 2024. URL <https://openreview.net/forum?id=GI5j6OMTju>.
- 553
- 554 Andrey Gromov, Kushal Tirumala, Hassan Shapourian, Paolo Glorioso, and Daniel A Roberts. The  
555 unreasonable ineffectiveness of the deeper layers. *arXiv preprint arXiv:2403.17887*, 2024.
- 556
- 557 Michael Hanna, Ollie Liu, and Alexandre Variengien. How does GPT-2 compute greater-than?:  
558 Interpreting mathematical abilities in a pre-trained language model. In *Thirty-seventh Confer-  
559 ence on Neural Information Processing Systems*, 2023. URL [https://openreview.net/  
560 forum?id=p4PckNQR8k](https://openreview.net/forum?id=p4PckNQR8k).
- 561 Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. Sparse  
562 autoencoders find highly interpretable features in language models. In *The Twelfth International  
563 Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?  
564 id=F76bwRSLeK](https://openreview.net/forum?id=F76bwRSLeK).
- 565 Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. What does BERT learn about the structure  
566 of language? In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the  
567 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3651–3657, Florence,  
568 Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1356. URL  
569 <https://aclanthology.org/P19-1356>.
- 570
- 571 Adam Karvonen, Benjamin Wright, Can Rager, Rico Angell, Jannik Brinkmann, Logan Riggs  
572 Smith, Claudio Mayrunk Verdun, David Bau, and Samuel Marks. Measuring progress in dictio-  
573 nary learning for language model interpretability with board game models. In *ICML 2024 Work-  
574 shop on Mechanistic Interpretability*, 2024. URL [https://openreview.net/forum?id=  
575 qzsDKwGJyB](https://openreview.net/forum?id=qzsDKwGJyB).
- 576
- 576 Connor Kissane, Ryan Krzyzanowski, Andrew Conmy, and Neel Nanda.  
577 SAEs (usually) transfer between base and chat models. [https:  
578 //www.alignmentforum.org/posts/fmwk6qxrPw8d4jvbd/  
579 saes-usually-transfer-between-base-and-chat-models](https://www.alignmentforum.org/posts/fmwk6qxrPw8d4jvbd/saes-usually-transfer-between-base-and-chat-models), July 2024. AI  
580 Alignment Forum.
- 581
- 581 Tom Lieberum, Senthoooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant  
582 Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse  
583 autoencoders everywhere all at once on gemma 2. *arXiv preprint arXiv:2408.05147*, 2024.
- 584
- 584 Aleksandar Makelov, Georg Lange, and Neel Nanda. Towards principled evaluations of sparse  
585 autoencoders for interpretability and control. In *ICLR 2024 Workshop on Secure and Trust-  
586 worthy Large Language Models*, 2024. URL [https://openreview.net/forum?id=  
587 MHIX9H8aYF](https://openreview.net/forum?id=MHIX9H8aYF).
- 588
- 588 Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller.  
589 Sparse feature circuits: Discovering and editing interpretable causal graphs in language models.  
590 *arXiv preprint arXiv:2403.19647*, 2024.  
591
- 592
- 592 Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous space word  
593 representations. In Lucy Vanderwende, Hal Daumé III, and Katrin Kirchhoff (eds.), *Proceedings*

- 594 of the 2013 Conference of the North American Chapter of the Association for Computational Lin-  
595 guistics: *Human Language Technologies*, pp. 746–751, Atlanta, Georgia, June 2013. Association  
596 for Computational Linguistics. URL <https://aclanthology.org/N13-1090>.
- 597  
598 Neel Nanda. Attribution patching: Activation patching at industrial scale, 2023.  
599 URL [https://www.neelnanda.io/mechanistic-interpretability/  
600 attribution-patching](https://www.neelnanda.io/mechanistic-interpretability/attribution-patching).
- 601 Neel Nanda, Andrew Lee, and Martin Wattenberg. Emergent linear representations in world models  
602 of self-supervised sequence models. In Yonatan Belinkov, Sophie Hao, Jaap Jumelet, Najoung  
603 Kim, Arya McCarthy, and Hosein Mohebbi (eds.), *Proceedings of the 6th BlackboxNLP Work-  
604 shop: Analyzing and Interpreting Neural Networks for NLP*, pp. 16–30, Singapore, December  
605 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.blackboxnlp-1.2. URL  
606 <https://aclanthology.org/2023.blackboxnlp-1.2>.
- 607 Frank Nielsen. *Hierarchical Clustering*, pp. 195–211. 02 2016. ISBN 978-3-319-21902-8. doi:  
608 10.1007/978-3-319-21903-5.8.
- 609  
610 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.  
611 Zoom in: An introduction to circuits. *Distill*, 2020. doi: 10.23915/distill.00024.001.  
612 <https://distill.pub/2020/circuits/zoom-in>.
- 613 Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry  
614 of large language models. In *Causal Representation Learning Workshop at NeurIPS 2023*, 2023.  
615 URL <https://openreview.net/forum?id=T0PoOJg8cK>.
- 616  
617 Judea Pearl. Direct and indirect effects. In *Probabilistic and causal inference: the works of Judea  
618 Pearl*, pp. 373–392. 2022.
- 619 Senthoran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János  
620 Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse  
621 autoencoders, 2024. URL <https://arxiv.org/abs/2407.14435>.
- 622  
623 Lee Sharkey, Dan Braun, and Beren Millidge. Taking the temperature of transformer circuits,  
624 2023. URL [https://www.alignmentforum.org/posts/z6QQJbtpkEAX3Aojj/  
625 interim-research-report-taking-features-out-of-superposition](https://www.alignmentforum.org/posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposition). Ac-  
626 cessed: 2024-08-18.
- 627 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In  
628 *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.
- 629  
630 Aaquib Syed, Can Rager, and Arthur Conmy. Attribution patching outperforms automated circuit  
631 discovery. In *The 7th BlackboxNLP Workshop*, 2024. URL [https://openreview.net/  
632 forum?id=RysbaxAnc6](https://openreview.net/forum?id=RysbaxAnc6).
- 633 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow,  
634 and Rob Fergus. Intriguing properties of neural networks, 2014. URL [https://arxiv.org/  
635 abs/1312.6199](https://arxiv.org/abs/1312.6199).
- 636 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya  
637 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard  
638 Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex  
639 Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, An-  
640 tonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo,  
641 Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric  
642 Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Hen-  
643 ryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski,  
644 Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu,  
645 Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee,  
646 Lucas Dixon, Machel Reid, Maciej Mikula, Mateo Wirth, Michael Sharman, Nikolai Chinaev,  
647 Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko  
Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo

- 648 Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree  
649 Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech  
650 Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh  
651 Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin  
652 Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah  
653 Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on  
654 gemini research and technology, 2024a. URL <https://arxiv.org/abs/2403.08295>.
- 655 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu-  
656 patiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Fer-  
657 ret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Char-  
658 line Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin,  
659 Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur,  
660 Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchi-  
661 son, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge,  
662 Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar,  
663 Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Wein-  
664 berger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang,  
665 Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin,  
666 Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen  
667 Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway,  
668 Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez,  
669 Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mo-  
670 hamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir  
671 Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leti-  
672 cia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Mar-  
673 tins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth,  
674 Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi,  
675 Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khat-  
676 wani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Os-  
677 car Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko  
678 Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana,  
679 Reza Ardeshtir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah  
680 Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth,  
681 Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Ko-  
682 cisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren  
683 Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao  
684 Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris  
685 Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine  
686 Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu,  
687 Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen  
688 Kenealy, Robert Dadashi, and Alek Andreev. Gemma 2: Improving open language models at a  
689 practical size, 2024b. URL <https://arxiv.org/abs/2408.00118>.
- 690 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen,  
691 Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L  
692 Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers,  
693 Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan.  
694 Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Trans-  
695 former Circuits Thread*, 2024. URL [https://transformer-circuits.pub/2024/  
696 scaling-monosemanticity/index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).
- 697 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
698 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von  
699 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-  
700 vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,  
701 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/  
file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt.  
Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In

702        *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=NpsVSN6o4ul>.  
703  
704  
705        Larry Wasserman. *All of Statistics: A Concise Course in Statistical Inference*. Springer Publishing  
706        Company, Incorporated, 2010. ISBN 1441923225.  
707        Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In  
708        *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12,*  
709        *2014, Proceedings, Part I 13*, pp. 818–833. Springer, 2014.  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A HYPERPARAMETERS

Table 1: Pythia-160M model specifics

Config	Value
Layers ( $L$ )	12
Model dimension ( $d_{\text{model}}$ )	768
Heads ( $H$ )	12
Non-Embedding params	85,056,000
Equivalent models	GPT-Neo OPT-125M

Table 2: Training and fine-tuning hyperparameters

Hyperparameter	Value
c	8
$\lambda$	3e-4
Hook name	resid-post
Batch size	4096
Adam ( $\beta_1, \beta_2$ )	(0, 0.999)
Context size	1024
lr	7e-4
lr scheduler	constant
lr decay steps	last 20% of the training steps
ll warm-up steps	5% of the training steps
# tokens (Train)	1B
Checkpoint freq	200M
Decoder weights initialization	Zeroes
Activation function	JumpReLU
Decoder column normalization	Yes
Activation normalization	Mean squared $\ell_2$ norm of one during SAE training
FP precision	32
Prepend BOS token	No

## B ADDITIONAL ANGULAR DISTANCES AND LAYERS GROUPS

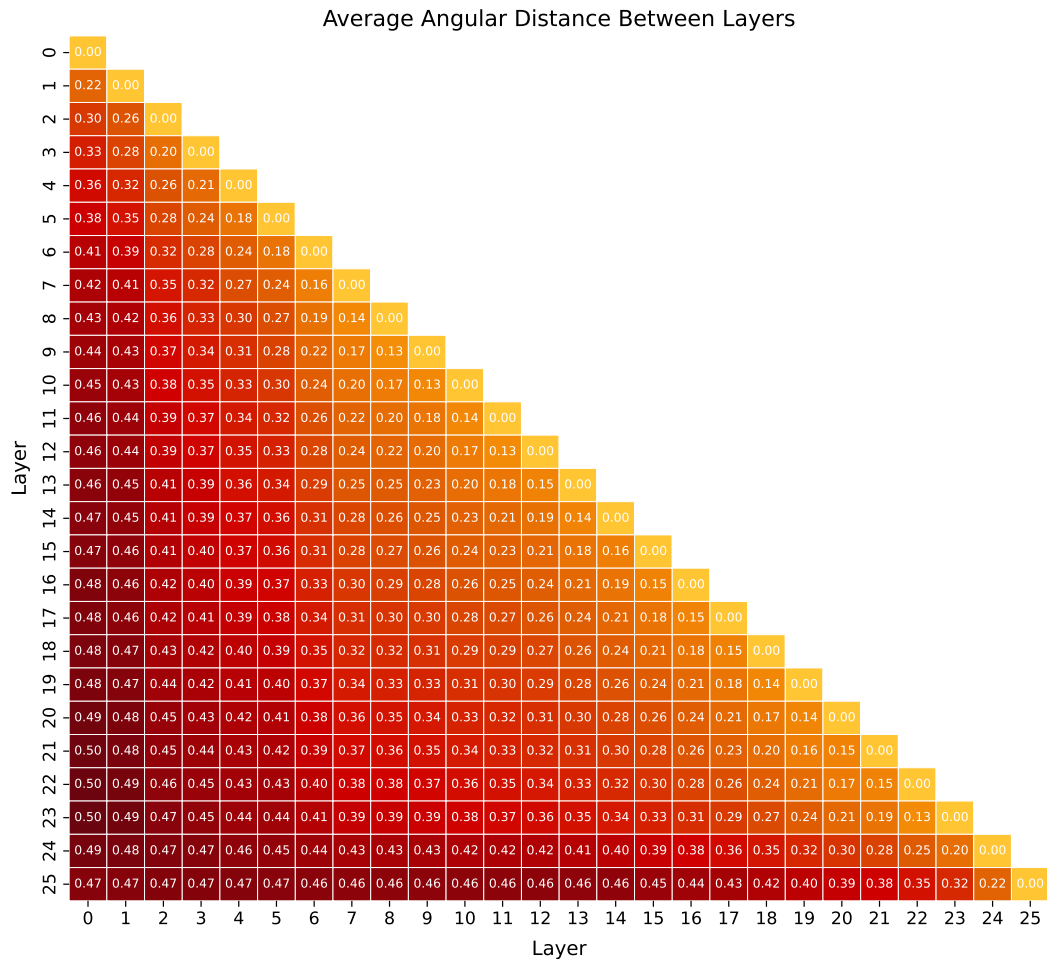


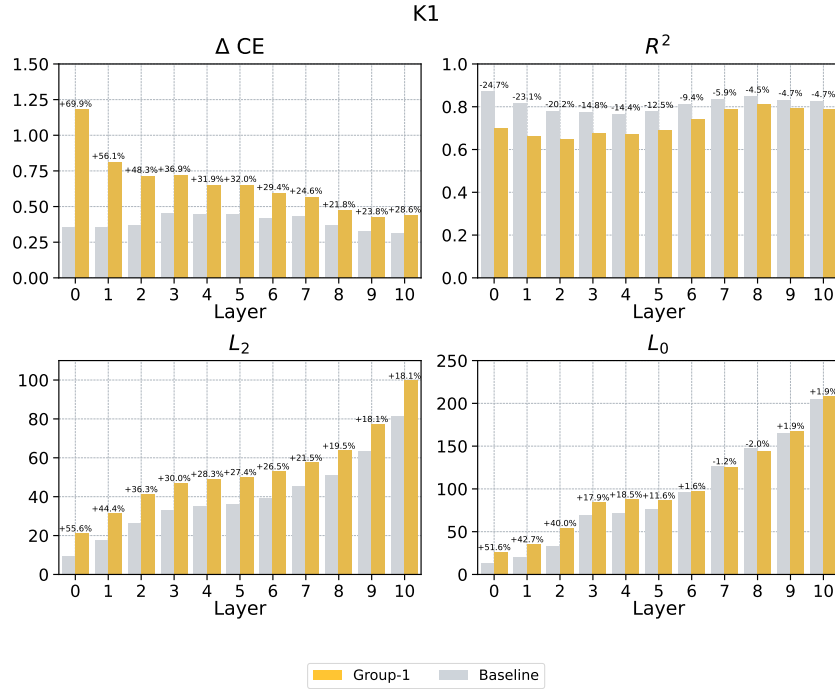
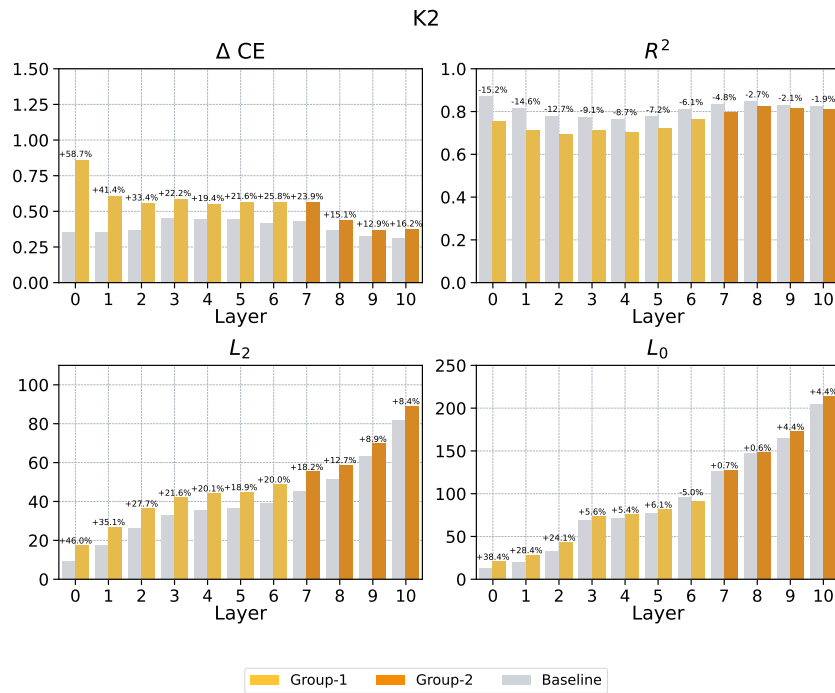
Figure 7: Average angular distance between all layers of Gemma-2 2b, as defined in Section 3.1. The angular distances are bounded in  $[0, 1]$ , where an angular distance equal to 0 means equal activations, 0.5 means activations are perpendicular and an angular distance of 1 means that the activations point in opposite directions.



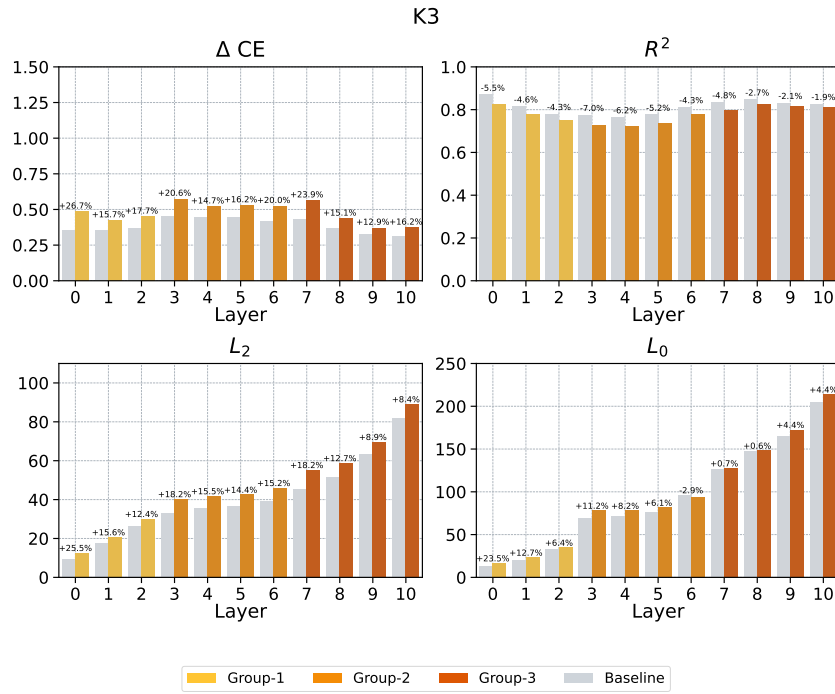
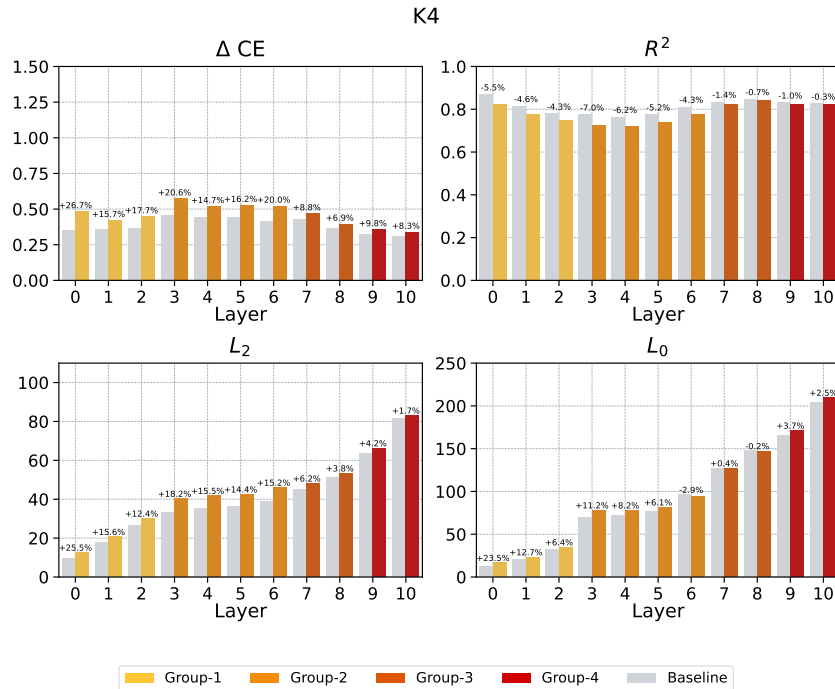




## C DETAILED PER-LAYER RECONSTRUCTION AND SPARSITY PLOTS

Figure 9: Per-layer  $\Delta CE$ ,  $R^2$ ,  $L_2$  and  $L_0$  with  $k = 1$ .Figure 10: Per-layer  $\Delta CE$ ,  $R^2$ ,  $L_2$  and  $L_0$  with  $k = 2$ .

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

Figure 11: Per-layer  $\Delta CE$ ,  $R^2$ ,  $L_2$  and  $L_0$  with  $k = 3$ .Figure 12: Per-layer  $\Delta CE$ ,  $R^2$ ,  $L_2$  and  $L_0$  with  $k = 4$ .

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

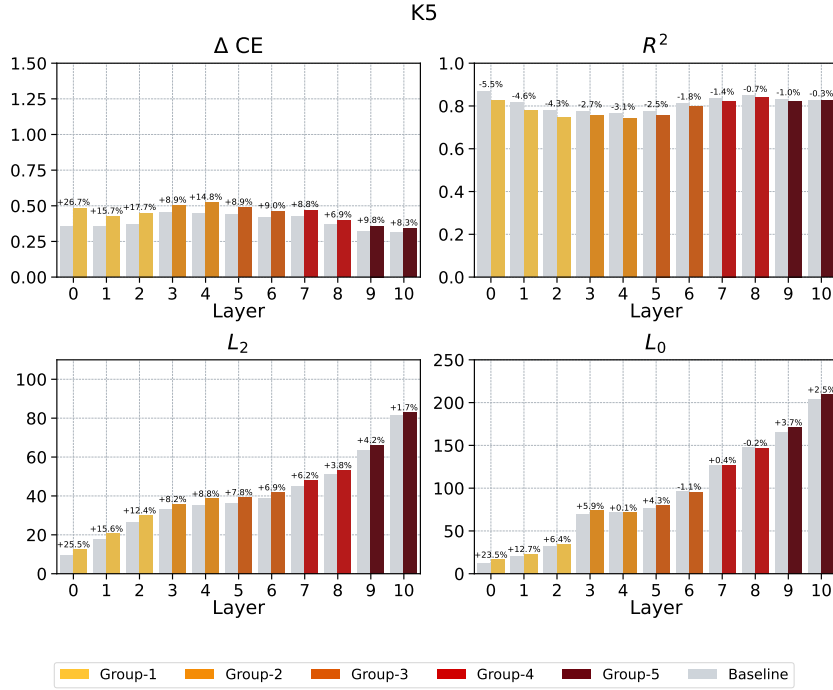


Figure 13: Per-layer  $\Delta CE$ ,  $R^2$ ,  $L_2$  and  $L_0$  with  $k = 5$ .

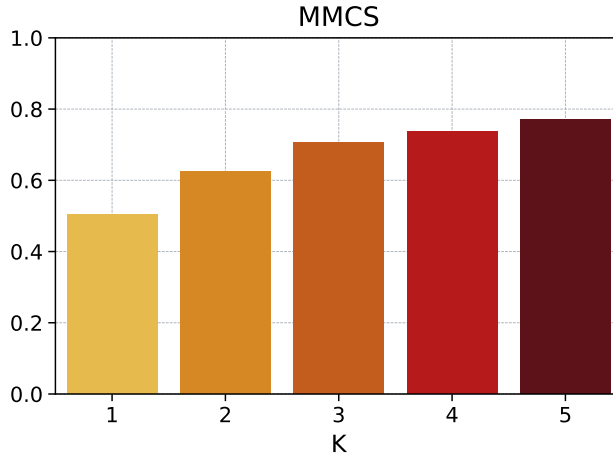


Figure 14: Per-layer  $\Delta CE$ ,  $R^2$ ,  $L_2$  and  $L_0$  with  $k = 5$ .

## D DETAILED PER-LAYER MMCS PLOTS

The Maximum Mean Cosine Similarity (MMCS) Score is defined as:

$$\text{MMCS} = \frac{1}{d_{\text{sae}}} \sum_{\mathbf{u}} \max_{\mathbf{v}} \text{CosSim}(\mathbf{u}, \mathbf{v}) \quad (7)$$

where  $\mathbf{u}$  and  $\mathbf{v}$  are the columns of the  $\text{SAE}_{j_k}$  and  $\text{SAE}_i$  decoder matrices, respectively. Figure 14 shows the average MMCS, where the average is computed for a given  $k$  by calculating the MMCS between  $\text{SAE}_{j_k}$  and  $\text{SAE}_s$  for every  $1 \leq j \leq k$  and  $s \in [j_k]$ , then dividing by  $L - 1$ .

Figure 15 shows the per-layer MMCS of every  $\text{SAE}_{j_k}$  for every  $k$ .

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

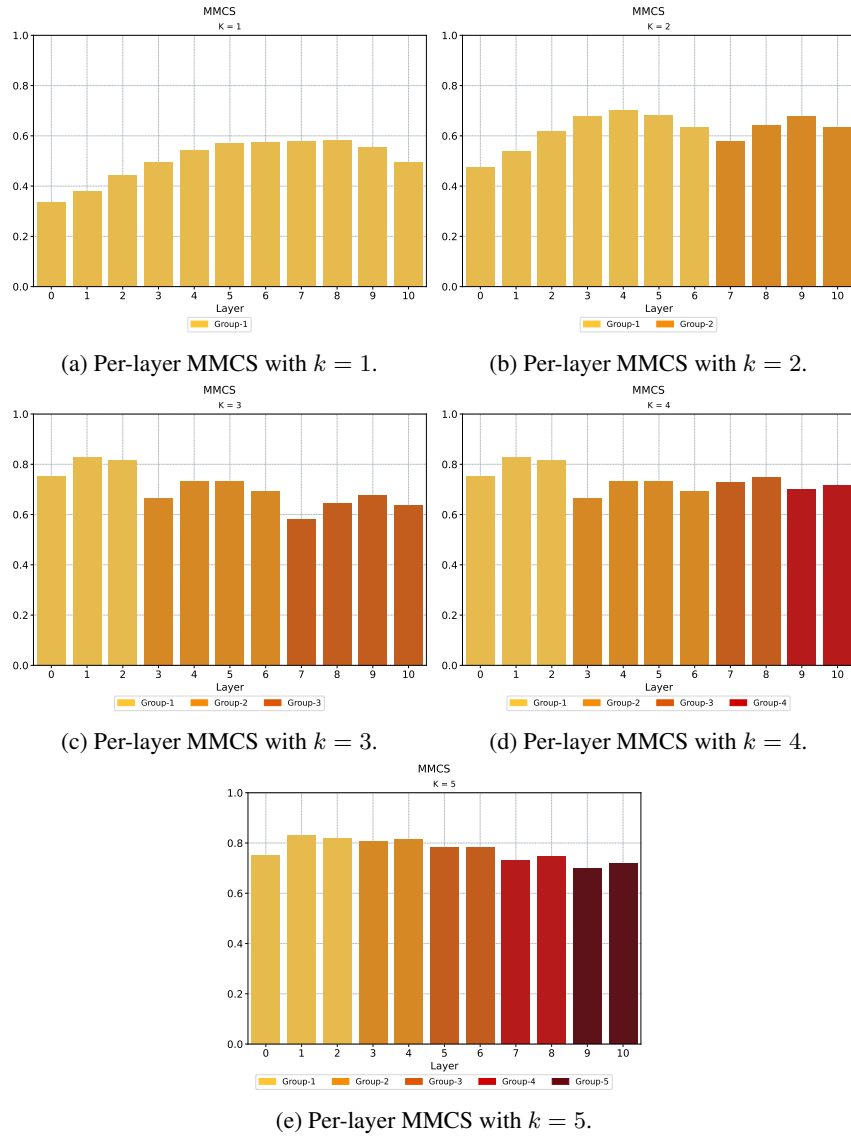


Figure 15: Per-layer MMCS for every  $k$ .

## E APPROXIMATE INDIRECT EFFECTS

In Equation 6 we reported the Indirect Effect (IE) (Pearl, 2022), which measures the importance of a feature with respect to a generic downstream task  $\mathcal{T}$ . To reduce the computational burden of estimating the IE with a single forward pass per feature, we employed two approximate methods: Attribution Patching (AtP) (Nanda, 2023; Syed et al., 2024) and Integrated Gradients (IG) (Sundararajan et al., 2017).

AtP (Nanda, 2023; Syed et al., 2024) employs a first-order Taylor expansion

$$\hat{\text{IE}}_{\text{AtP}}(m; \mathbf{f}_i; x_{\text{clean}}, x_{\text{corrupted}}) = \nabla m|_{\mathbf{f}_i = \mathbf{f}_{i;\text{clean}}}(\mathbf{f}_{i;\text{corrupted}} - \mathbf{f}_{i;\text{clean}}) \quad (8)$$

which estimates Equation 6 for every  $\mathbf{f}_i$  in two forward passes and a single backward pass.

Integrated Gradients (Sundararajan et al., 2017) is a more expensive but more accurate approximation of Equation 6

$$\hat{\text{IE}}_{\text{IG}}(m; \mathbf{f}_i; x_{\text{clean}}, x_{\text{corrupted}}) = \left( \sum_{\alpha \in \Gamma} \nabla m|_{\alpha \mathbf{f}_{i;\text{clean}} + (1-\alpha)\mathbf{f}_{i;\text{corrupted}}} \right) (\mathbf{f}_{i;\text{corrupted}} - \mathbf{f}_{i;\text{clean}}) \quad (9)$$

where  $\alpha$  ranges in an equally spaced set  $\Gamma = \{0, \frac{1}{N}, \dots, \frac{N-1}{N}\}$ . In our experiments we have set  $N = 10$ .

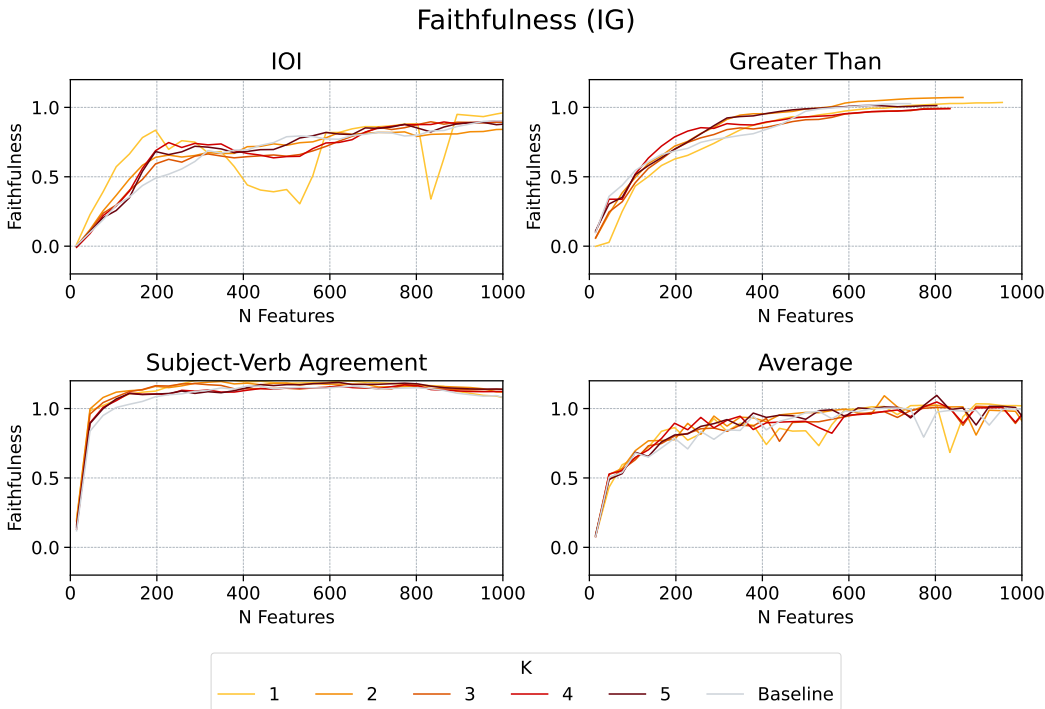


Figure 16: Average faithfulness (Marks et al., 2024) for every downstream task (IOI, Greater Than, Subject-Verb Agreement). The “Baseline” average is computed considering the performance obtained by  $\text{SAE}_i, \forall i = 0, \dots, 10$ . The “Average” plot depicts the average over the three downstream tasks.

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

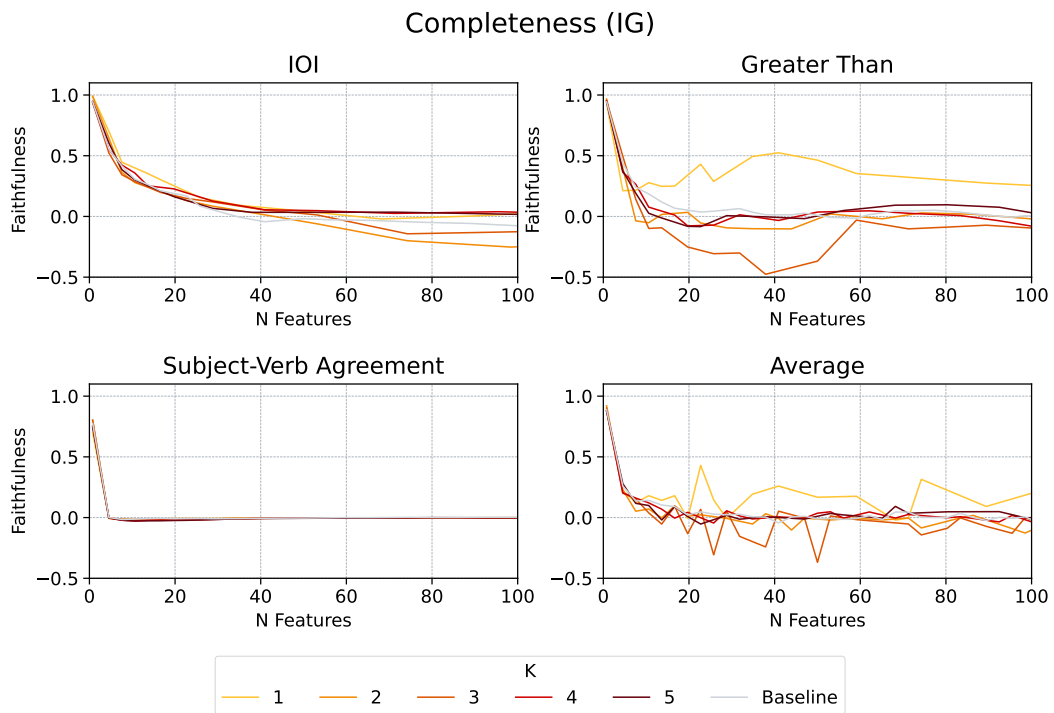


Figure 17: Average completeness for every downstream task (IOI, Greater Than, Subject-Verb Agreement). The “Baseline” average is computed considering the performance obtained by  $SAE_i$ ,  $\forall i = 0, \dots, 10$ . The “Average” plot depicts the average performance over the three downstream tasks.



## F HUMAN INTERPRETABILITY SCORES

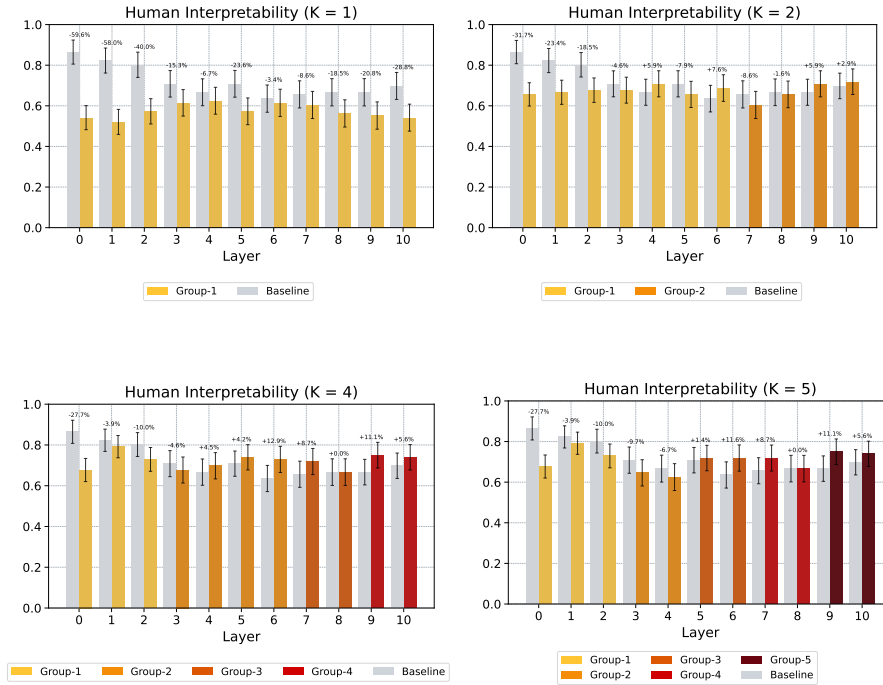


Figure 18: Human Interpretability Scores for  $k \in \{1, 2, 4, 5\}$ . The differences in the interpretability scores of features learned by the  $SAE_{j_k}$  and the baseline  $SAE_i$  are not statistically significant different from zero for most of the layers. Error bars shows one standard deviations of the scores differences, modeled as a Binomial distribution (Wasserman, 2010)

1350 G SCALING GROUP SAES

1351

1352 Results on Gemma-2 2b will be added here.

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403