

A IS FOR ABSORPTION: STUDYING FEATURE SPLITTING AND ABSORPTION IN SPARSE AUTOENCODERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Sparse Autoencoders (SAEs) have emerged as a promising approach to decompose the activations of Large Language Models (LLMs) into human-interpretable latents. In this paper, we pose two questions. First, to what extent do SAEs extract monosemantic and interpretable latents? Second, to what extent does varying the sparsity or the size of the SAE affect monosemanticity / interpretability? By investigating these questions in the context of a simple first-letter identification task where we have complete access to ground truth labels for all tokens in the vocabulary, we are able to provide more detail than prior investigations. Critically, we identify a problematic form of feature-splitting we call “feature absorption” where seemingly monosemantic latents fail to fire in cases where they clearly should. Our investigation suggests that varying SAE size or sparsity is insufficient to solve this issue, and that there are deeper conceptual issues in need of resolution. We release a feature absorption explorer at <https://feature-absorption.streamlit.app>.

1 INTRODUCTION

Large Language Models (LLMs) have achieved remarkable performance across a wide range of tasks, yet our understanding of their internal mechanisms lags behind their capabilities. This gap between performance and interpretability raises concerns about the “black box” nature of these models (Rudin, 2019). The field of mechanistic interpretability aims to address this issue by reverse-engineering the internal algorithms of neural networks and performing causal analysis on them (Olah et al., 2020).

One recent promising approach in this field is the use of Sparse Autoencoders (SAEs), which have shown potential in decomposing the dense, polysemantic activations of LLMs into more “interpretable” latent features (Cunningham et al., 2024; Bricken et al., 2023) using sparse dictionary learning (Olshausen & Field, 1997). SAE neurons (hereafter called “latents”) ¹ are said to be interpretable if they appear to detect some property of the input (which we refer to as a “feature”) and classify that feature with high precision / recall (Bricken et al., 2023).

¹We use *latents* to prevent overloading the term *feature*, which we reserve for human-interpretable concepts the SAE may capture. This breaks from earlier usage which used *feature* for both (Elhage et al., 2022), but aligns with the terminology in (Lieberum et al., 2024) and makes the distinction more clear.

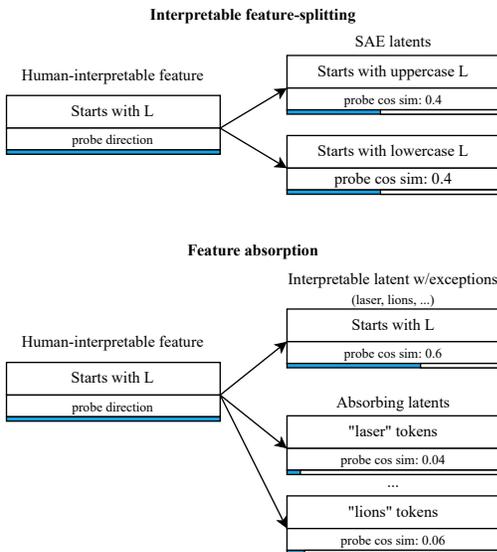


Figure 1: In feature absorption, an SAE latent appears interpretable, but has arbitrary exceptions where different latents “absorb” the feature direction and activate instead. The absorbing latents are frequently, though not always, token-aligned.

054 However, despite these theoretical advantages, most existing work on SAE interpretability mainly
 055 studies max activating examples (Cunningham et al., 2024), which may be misleading. There are
 056 more rigorous works which only measure precision (Bricken et al., 2023; Templeton et al., 2024;
 057 Kissane et al., 2024). Recent work has briefly explored recall and found it to be worse than naively
 058 expected, but this remains poorly understood (Olah et al., 2024b). We build on this work by eval-
 059 uating precision / recall on a large number of SAEs, and offer a partial explanation for lower-than-
 060 expected recall of SAE latents in the form of “feature absorption”.

061 Our key contributions in this investigation include the following:

- 062 1. We identify numerous SAE latents appearing to classify first-letter features. We calculate
 063 their precision / recall on the first letter identification tasks (as a proxy for monosemanticity
 064 / interpretability) and find they significantly underperform linear probes.
- 065 2. We find that latents which seem outwardly to classify the same feature can have vastly
 066 different precision / recall and that this tradeoff is mediated by various factors, mainly
 067 sparsity and width of the SAE.
- 068 3. Most importantly, we identify and quantify a variant of feature-splitting we call “feature
 069 absorption”, where an SAE latent appears to track a human-interpretable concept, but it
 070 fails to activate on seemingly arbitrary tokens. Instead, a different latent activates and
 071 contributes a portion of the probe direction, “absorbing” the feature. This is described in
 072 Figure 1.

073 We believe that feature absorption poses an obstacle to the practical application of SAEs since it
 074 suggests SAE latents may be inherently unreliable classifiers. This is particularly important if we
 075 seek to use them in safety applications where we need confidence that latents are fully tracking
 076 behaviours, such as bias or deceptive behavior. Furthermore, techniques which seek to describe
 077 circuits in terms of a sparse combination of latents with also be more difficult in the context of
 078 feature absorption (Marks et al., 2024).

081 2 BACKGROUND

082
 083 **Linear probing.** A linear probe is a simple linear classifier trained on the hidden activations of a
 084 neural network, typically using logistic regression (LR) (Alain & Bengio, 2017).

085
 086 **K-sparse probing.** A k-sparse probe (Gurnee et al., 2023) is a linear probe trained on a sparse
 087 subset of k neurons or SAE latents. Training a k-sparse probe first requires selecting the k best
 088 neurons or SAE latents that in-aggregate act as a good classifier, and then training a standard linear
 089 probe on just those k neurons or latents.

090 Gurnee et al. (2023) proposed several methods of estimating the best k neurons or features to pick,
 091 one of which involves first training a LR probe with a L1 loss term, and selecting the k largest
 092 elements by probe weight. When we refer to k-sparse probing in this work, we use this method of
 093 selecting k features.

094
 095 **Sparse autoencoders.** An SAE consists of an encoder, W_{enc} , a decoder, W_{dec} , and correspond-
 096 ing biases b_{enc} and b_{dec} . The SAE has a nonlinearity, σ , typically a ReLU (or variants such as
 097 JumpReLU (Rajamanoharan et al., 2024; Lieberum et al., 2024)). Given input activation, a , the
 098 SAE computes a hidden representation, f , and reconstruction, \hat{a} :

$$100 \quad f = \sigma(W_{enc}a + b_{enc}) \quad (1)$$

$$101 \quad \hat{a} = W_{dec}f + b_{dec} \quad (2)$$

102
 103 SAEs attempt to reconstruct input activations by projecting into an overcomplete basis using a
 104 sparsity-inducing loss term (typically $L1$ loss), or a certain number of non-zero features ($L0$) on
 105 the hidden activations. SAEs learn feature decompositions in an unsupervised manner, and while
 106 the sparsity penalty is meant to encourage monosemantic features, it is often hard to judge if the fea-
 107 tures learned are interpretable or to say with certainty that the features the SAE learned are faithful
 to the computation performed by the underlying LLM.

SAE feature ablation. We often want to understand how an SAE latent causally influences a downstream output. In an ablation study, the latent in question is removed from the computation graph of the model to see the effect this has on a downstream metric. A negative ablation effect means removing the SAE latent would lower the metric.

We follow the work of Marks et al. (2024) and provide the procedure in Algorithm 1 below. We also make use of the integrated-gradients (IG) approximation (Sundararajan et al., 2017) to improve the speed of running multiple ablation experiments.

Algorithm 1 SAE Latent Ablation

```

1: Insert SAE in model computation path, including error term
2: Define a scalar metric on the model’s output distribution (e.g. difference between token logits)
3: Calculate baseline metric value for a test prompt
4: for each token of interest do
5:   for each SAE latent do
6:     Set the SAE latent activation to 0
7:     Recalculate the metric
8:     Compute ablation effect as (baseline metric - new metric)
9:     Reset the SAE latents to its original value
10:  end for
11: end for

```

3 EXPERIMENTAL SETUP

Our experiments focused on predicting the first-letter of a single token containing characters from the English alphabet (a-z, A-Z) and an optional leading space. We use in-context learning (ICL) prompts to elicit knowledge from the model, using templates of the form:

```
{token} has the first letter: {capitalized_first_letter}
```

An example of an ICL prompt consisting of 2 in-context examples is shown below. The model should output the `_D` token:

```

tartan has the first letter: T
mirth has the first letter: M
dog has the first letter:

```

In the above prompt, we extract residual stream activations at the `_dog` token index. These activations are used both for LR probe training and for applying SAEs. We use a train/test split of 80% / 20%, and evaluate only on the test set, including when running experiments on SAEs. When applying SAEs, we include the SAE error term (Marks et al., 2024) to avoid changing model output.

To determine the causal effect of SAE latents on the first-letter identification task we conduct ablation studies. We use a metric consisting of the logit of the correct letter minus the mean logit of all incorrect letters. This measures the propensity of the model to choose the correct starting letter as opposed to other letters. Formally, our metric m is defined below, where g refers to the final token logits, L is the set of uppercase letters, and y is the uppercase letter that is the correct starting letter:

$$m = g[y] - \frac{1}{|L| - 1} \sum_{l \in \{L \setminus y\}} g[l]$$

We discuss this metric and alternative formulations further in Appendix A.8.

To determine how well multiple features perform as a classifier when used together, we use k-sparse probing, increasing the value of k from 1 to 15. We train a LR probe using a L1 loss term with coefficient 0.01, and select the top k features by magnitude.

We use the base Gemma-2-2B model for most of our studies, along with the full set of Gemma Scope residual stream SAEs of width 16k and 65k released by Deepmind (Lieberum et al., 2024).

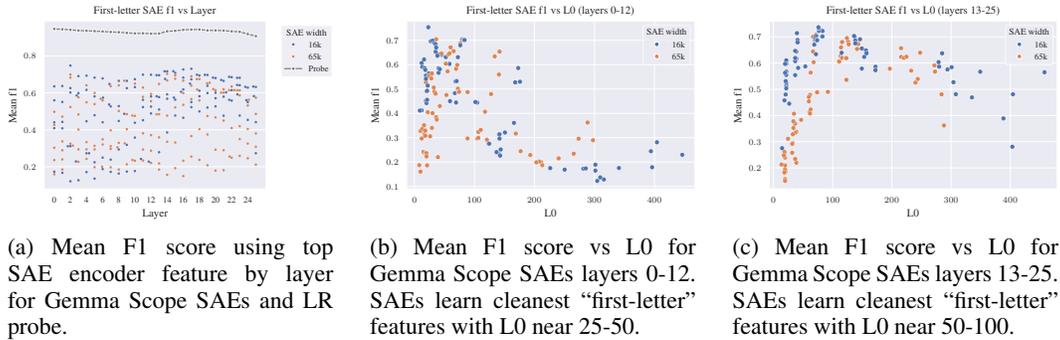


Figure 2: Comparison of F1 scores for first-letter classification tasks

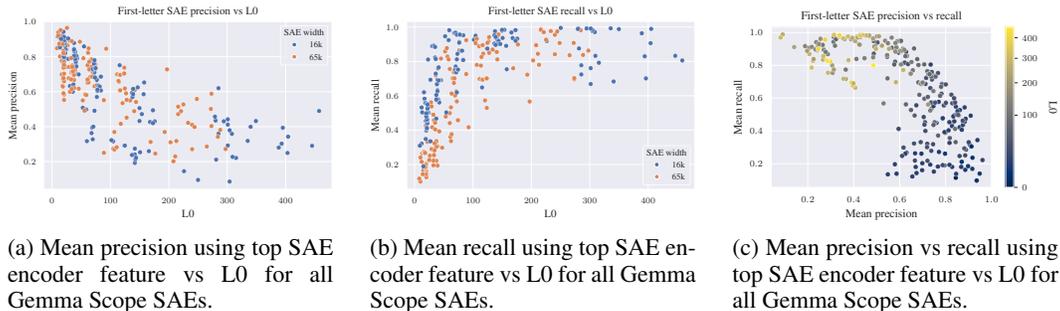


Figure 3: Precision and recall vs L0 for first-letter classification tasks

We also evaluate absorption on our own SAEs trained on Qwen2 0.5B (Yang et al., 2024) and Llama 3.2 1B (Dubey et al., 2024).

4 RESULTS

Our results are divided into three sections. First, we compare the performance of linear probes with SAE latents on recovering first-character information from model activations, showing that despite appearing to track first letter features, a wide variety of precision / recall is achieved. Second, we motivate our definition of feature absorption with a case-study, emphasizing how an absorbing feature can unexpectedly causally mediate first letter information whilst the first-letter latent (unexpectedly) fails to fire. Finally, we attempt to quantify feature splitting and feature absorption, showing that tuning of hyper-parameters may partially assist but not fully alleviate feature absorption.

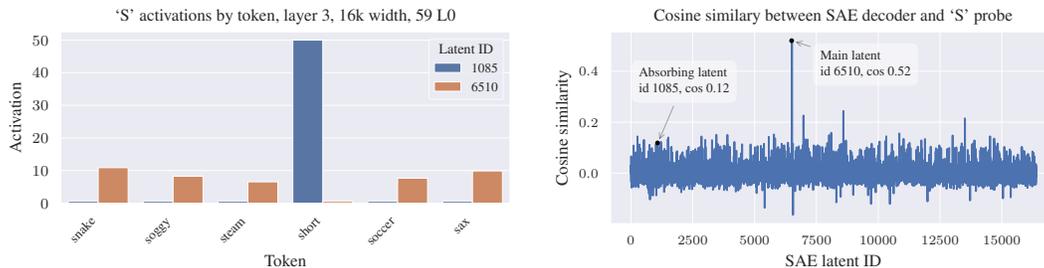
4.1 DO SAES LEARN LATENTS THAT TRACK FIRST LETTER INFORMATION?

We compare the performance of LR probes with the performance of the SAE latent whose encoder direction has highest cosine similarity with the probe, resulting in 26 “first-letter” latents. We observed that for each probe, there was clearly one or at most a couple of outlier SAE latents with high probe cosine similarity. Full plots of cosine similarity vs letter are shown in Appendix A.5.

We also experimented with using $k=1$ sparse probing to identify SAE latents (Gurnee et al., 2023), and find this gives similar results. Further comparison of $k=1$ sparse probing and encoder cosine similarity is explored in Appendix A.4.

We observe wide variance in the performance of Gemma Scope SAEs at the first-letter identification task, but no SAE matches LR probe performance. We show the mean F1 score by layer as well as the F1 score of the LR probe in Figure 2a. We further investigate the F1 score of these SAE encoder latents as a function of L0 and SAE width in Figures 2b and 2c.

Whether or not an SAE learns a clear “first-letter” latent for each letter is highly dependent on L0, with low L0 SAEs tending to learn high-precision low-recall latents, and high L0 SAEs learning



(a) Layer 3, L0=59 SAE feature activations for tokens that start with “S”. The core “starts with S” feature, 6510, fails to activate on the token `_short`. The “short”-token aligned feature 1085 activates instead.

(b) Cosine similarity between layer 3, L0=59 SAE decoder and the “starts with S” probe. The main “Starts with S” latent, 6510, is clearly visible and highly probe-aligned.

Figure 4: SAE activations and cosine similarity for “starts with S” features.

low-precision high-recall latents (Figure 3). We caution drawing conclusions about an “optimal” L0 from these plots, as we find further variance when broken-down by letter, shown in Appendix A.5.

4.2 WHY DO SAE LATENTS UNDERPERFORM?

The Gemma Scope layer 3, 16k width, 59 L0 SAE has a latent, 6510, which appears to act as a classifier for “starts with S”, achieving an F1 of 0.81. However, this latent fails to activate on some tokens the probe can classify, and which the model can spell, such as the token `_short`.

Figure 4a shows a sample prompt containing a series of tokens that start with “S”, and the activations of top SAE latents by ablation score for these tokens. The main “starts with S” latent, 6510, activates on all these tokens except `_short`. This SAE also has a token-aligned latent, 1085, which activates on variants of the word “short” (“short”, “SHORT”, etc...). The Neuronpedia dashboard (Lin & Bloom, 2023) for feature 1085 is shown in Appendix A.11. For the token `_short`, the main “starts with S” latent does not activate but the “short” latent activates instead.

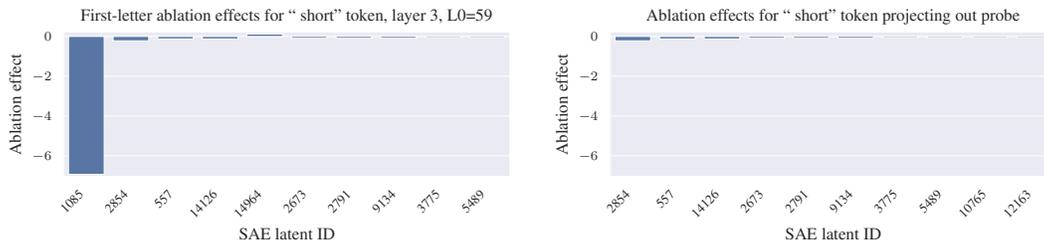
Latent 1085 has a cosine similarity with the “starts with S” probe of 0.12, indicating it contains a component of the “starts with S” direction, although much smaller than the main “starts with S” latent. Cosine similarity of the SAE decoder with the “starts with S” LR probe is shown in Figure 4b. Interestingly, despite latent 1085 having only about 1/5 the cosine similarity with the probe as the main latent 6510, we see it activates with about 5 times the magnitude of latent 6510 on the `_short` token, thus contributing a similar amount of the “starts with S” probe direction to the residual stream.

We conduct an ablation experiment on the `_short` token, shown in Figure 5a, and see that latent 1085 has a dramatically larger ablation effect compared with all other SAE features. This suggests latent 1085 is causally responsible for the model knowing that `_short` starts with S.

Is it possible that the probe projection is not the causally important component of feature 1085? We conduct another ablation experiment, except now we remove the probe direction from feature 1085 via projection before ablation. The results of this ablation experiment are shown in Figure 5b. After removing the probe component from feature 1085, it no longer has a significant ablation effect. Thus we know the probe projection of feature 1085 is responsible for model behavior.

These experiments show the “starts with S” feature has been “absorbed” by the token-aligned latent 1085, likely along with other semantic concepts related to the word “short”. After observing that the main “starts with S” latent 6510 activates on most tokens that begin with “S”, it may be tempting to conclude this latent tracks the interpretable feature of beginning with the letter “S”. However, this latent quietly fails to activate on the `_short` token, leading us to a false sense of understanding.

We call this phenomenon **feature absorption**. In feature absorption a seemingly interpretable SAE latent fails to activate on arbitrary positive examples, and instead the feature is “absorbed” into approximately token-aligned latents.



(a) Ablation effect for `_short` token, indicating that feature 1085, is responsible for the “starts with S” concept for the `_short` token. The main “starts with S” latent, 6510, does not activate on the `_short` token.

(b) Ablation effect for `_short` token after removing the probe direction from latent 1085 via projection. Latent 1085 no longer appears in the plot, indicating the strong ablation effect in Figure 5a is due to its component along the probe direction.

Figure 5: Ablation effects on `_short` token before and after projecting out the probe direction

Latent 7112	Latent 7657
žda se naplaćuje naknada	LC, an aluminum boat
. E. Søli, 20	as LIFT and LF-Net. Once
a></code>	latter’s sister Louise, who in

Table 1: Sample max activating examples for latents 7112 and 7657 for Gemma Scope 16k, layer 0, 105 L0 from Neuronpedia. The token where the SAE feature activates is highlighted in yellow. Latent 7112 appears to be a lowercase “L” starting-letter latent, and latent 7657 appears to be a corresponding uppercase “L” latent.

Feature absorption is likely a logical consequence of SAE sparsity loss. If a dense and sparse feature co-occur, absorbing the dense feature into a latent tracking the sparse feature will increase sparsity.

4.3 MEASURING FEATURE SPLITTING AND FEATURE ABSORPTION

Feature splitting A key phenomenon identified from previous studies of SAEs is feature-splitting (Bricken et al., 2023), where a feature represented in a single latent in a smaller SAE can split into two or more latents in a larger SAE. During our experiments, we found strong evidence of feature-splitting in the Gemma Scope SAEs.

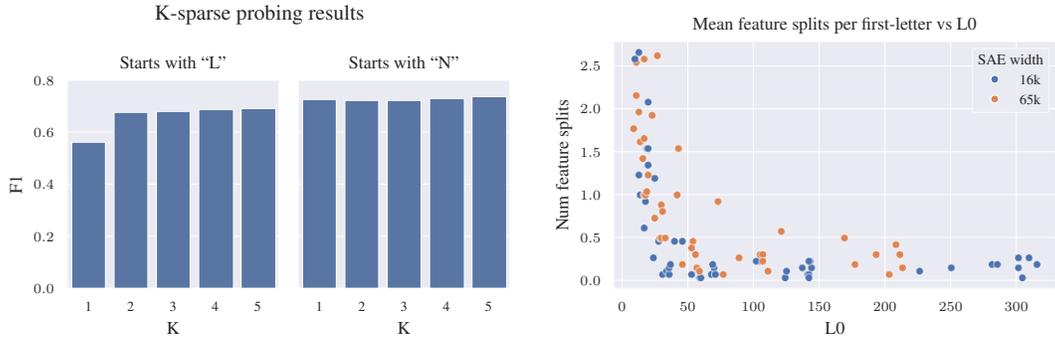
For instance, in the layer 0, 16k width, 105 L0 SAE, we find two encoder latents (id:7112 and id:7657²) which align with the “L” starting letter probe. Inspecting max activating examples, we see latent 7112 activates on tokens starting with lowercase “l”, while 7657 activates on tokens starting with uppercase “L”. Some activating examples for these latents are shown in Table 1.

Feature splitting like this is not necessarily problematic for interpretability efforts since the split features are still easily identifiable, and depending on the context it may be more useful to have either a single “starts with L” latent or a pair of “starts with uppercase / lowercase L” latents.

We measure feature splitting using k-sparse probing (Gurnee et al., 2023) on SAE activations. If increasing the k-sparse probe from k to $k + 1$ causes a significant increase in probe F1 score, then the additional SAE latent provides a meaningful signal, and the combination of these $k + 1$ latents is likely a feature split. In the example of the uppercase “L” and lowercase “l” split, a k-sparse probe with $k = 2$ trained on both these features should predict “starts with letter L” much better than either feature on its own. Figure 6a shows F1 vs K for letters “L” and “N”. The “L” k-sparse probe shows a significant jump in F1 score moving from $k=1$ to $k=2$ corresponding to feature splitting, while the F1 score for the “N” k-sparse probe is relatively constant.

We detect feature splitting by measuring whether increasing k by one causes a jump in F1 score by more than threshold τ . We set $\tau = 0.03$ after manually inspecting latents with various thresholds. Figure 6b shows feature splitting vs L0 for all 16k and 65k width Gemma Scope SAEs.

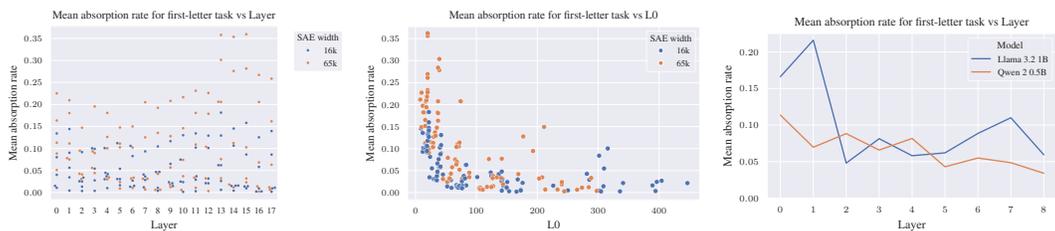
²<https://www.neuronpedia.org/list/cm0h1n2mt00019jdk274owq9e>



(a) K-sparse probing results for letters “L” and “N”, layer 0, 16k width, 105 L0. “L” shows a significant improvement in F1 between $k=1$ and $k=2$ corresponding to feature splitting.

(b) Mean number of feature splits per letter on the first-letter spelling task by L0. Feature splitting occurs more frequently with higher sparsity.

Figure 6: Feature splitting



(a) Mean feature absorption rate vs layer on the first-letter task, Gemma Scope 16k and 65k SAEs. We do not see an obvious pattern in absorption rates by layer.

(b) Mean feature absorption rate vs L0 on first-letter task, Gemma Scope 16k and 65k SAEs. Wider and more sparse SAEs demonstrate higher rates of absorption.

(c) Mean feature absorption rate vs layer on the first-letter task on Llama 3.2 1B and Qwen 2 0.5B models, standard L1 loss SAE architecture, layers 0-8.

Figure 7: Feature absorption rates

Feature absorption The single latent or a set of traditional feature split latents that seem to act as a classifier for a human-interpretable feature like “starts with S” fail to fire in a seemingly arbitrary number of cases. What fires instead are approximately token-aligned latents with small but positive alignment with the LR probe. We say these latents are absorbing the feature.

We quantify the extent to which feature absorption occurs with the metric **feature absorption rate**. We first find k feature splits for a first-letter feature using a k -sparse probe. We then find false-negative tokens that all k feature-split SAE latents fail to activate on, but which the LR probe correctly classifies, and run an integrated-gradients ablation experiment on those tokens. The ablation effect finds the most causally important SAE latents for the spelling of that token. If the SAE latent receiving the largest negative magnitude ablation effect has a cosine similarity with the LR probe above 0.025, and is at least 1.0 larger than the latent with the second highest ablation effect, we say that feature absorption has occurred. These thresholds were chosen from manual inspection of the data. We then calculate feature absorption rate as below:

$$\text{absorption_rate} = \frac{\text{num_absorptions}}{\text{lr_probe_true_positives}}$$

If there are more than 200 false negative per letter, we randomly sample 200 samples to estimate the number of absorptions. We see that absorption rate increases with higher sparsity and higher SAE width. Lower L0 likely pushes the SAE to absorb dense features like spelling information across multiple latents, thereby increasing feature sparsity. Feature absorption rate vs L0 for Gemma Scope SAEs layers 0-17 is shown in Figure 7b. Absorption rate by letter is shown in Appendix A.10. We also train our own set of standard L1 loss SAEs on the first 8 layers of Qwen2 0.5B (Yang et al., 2024) and Llama 3.2 1B Dubey et al. (2024). In Figure 7c we show that absorption occurs in these SAEs as well.

Our metric cannot capture absorption past layer 17 in Gemma 2 2B since we rely on ablation experiments to be certain the absorbed feature causally mediates model behavior. Past layer 17, attention has already moved the starting letter information from the source token into the final token position, so any ablations on the source token past layer 17 have little effect. This is a limitation of our absorption metric - we rely on ablation to be certain of the causal impact of absorbed features on model behavior, but this limits the layer depth our metric can be applied. We discuss this further in Appendix A.9.

Our absorption metric is not perfect, and is likely an under-estimate of the true level of feature absorption. We only consider absorption to have occurred if a single SAE latent has a much larger ablation effect than all other latents, and if the main SAE latents for a feature do not activate at all. Our metric will not capture multiple absorbing latents activating together, or the main latents activating but very weakly. Regardless, we feel our metric is a reasonable conservative baseline.

5 RELATED WORK

5.1 APPLICATIONS OF PROBES AND SAES FOR MODEL INTERPRETABILITY

Probing methods have often been used to extract interpretable information from language models. However, the existence of such a representation does not guarantee that the model relies on this representation in its computation graph (Elazar et al., 2021).

Prior work has shown that many human-interpretable concepts in LLM activations are represented as linear directions in activation space, known as the linear representation hypothesis (Elhage et al., 2022; Park et al., 2024). Li et al. (2023) used non-linear probes to recover board representations from a transformer trained on Othello scripts (“OthelloGPT”). However, Nanda (2023) later showed that linear representations were not only recoverable but also editable.

Recent work has focused on applying SAEs to extract human-interpretable explanations of model internals. Karvonen et al. (2024) investigated how SAEs represent board states of Chess and Othello, and introduce a coverage metric using ground-truth features to evaluate the quality of SAE latents. This is very similar to our technique for evaluating SAE latents on a known task.

Other work has noted poor recall / precision of SAE latents compared to known proxies (Olah et al., 2024b; Kissane et al., 2024; Templeton et al., 2024). We build on this work by evaluating a large number of Gemma Scope SAEs to demonstrate how precision / recall is mediated by sparsity, and also offer a possible explanation of low recall due to feature absorption.

5.2 CHARACTER-LEVEL INFORMATION IN LANGUAGE MODELS

The ability of LLMs to learn character-level information from ostensibly character-blind tokens has been studied by various scholars, though no clear mechanism has yet been established. Kaushal & Mahowald (2022) trained MLPs from the embedding layers of GPT-J as probes for each letter in the alphabet, again finding good performance that implied character-level information was represented, but did not look into the model internals to explain how these representations were being used. In a follow-up work, Watkins & Bloom (2023) demonstrated that even linear probes on the embedding layers perform comparably well to MLPs in extracting character-level information.

In contrast to the above approaches, we train various probes for each character on multiple layers, and compare with SAE latents for a variety of layers.

5.3 DECOMPOSING SAE LATENTS

The phenomenon of SAEs of different sides splitting a feature into various smaller latents was first described in Bricken et al. (2023), which noted that different SAE widths and sparsities induce latents of different granularity, with wider SAEs often learning more specific variants of features.

Bussmann et al. (2024) find that by training an SAE on the decoder of another SAE, a technique called Meta-SAEs, it is possible to break down a single SAE latent like “Einstein” into subcomponents like “German” and “Physicist” and “starts with E”. Meta-SAEs may provide a promising future research direction to overcome the feature absorption phenomenon we describe in this paper.

6 DISCUSSION

Interpretability of SAE latents In this work, we use a simple first-letter identification task to investigate whether SAEs extract monosemantic and interpretable features, and how this is affected by varying hyperparameters like SAE size, sparsity, or layer. We find that the SAE latents we investigated were not interpretable and that varying the sparsity or the size of the SAE did not meaningfully change this.

One may argue it is unreasonable to judge an SAE latent against a linear probe directly optimized to perform the same classification task and that it has been established that sometimes unsupervised methods can surprise us (Nanda, 2023). However, we argue that for a latent to be considered interpretable, its behavior should match what one would reasonably expect the latent to be doing after inspecting its activation patterns. In our experiments, we use SAE latents that do appear to perform first letter classification, and then evaluate how well they perform this task. We validate as well that these latents causally mediated model performance on the first-letter task in Appendix A.3. We are convinced that these latents should reasonably be considered “first-letter” latents and that their performance on the first-letter identification task is a valid measure of their interpretability.

Feature absorption In trying to understand why SAE latents fail to match the performance of LR probes, we identified a form of feature splitting we call “feature absorption”. Feature absorption may be particularly problematic for SAEs because it creates an interpretability illusion where we believe we have found an interpretable latent, but absorption induces lower recall by creating clear false negatives to the mainline interpretation of the latent. For example, we may believe we have found a SAE latent which tracks deceptive behavior in the model, but due to feature absorption, there may be many cases where that latent fails to fire. This lower recall poses problems for methods which rely on using SAEs to find sparse circuits (Marks et al., 2024), as the number of latents needed to characterize model behavior may be much larger than expected. We find that feature absorption happens even in high-recall latents, so this is not only a problem for low L0 SAEs and appears to be a more fundamental issue.

We hypothesize that feature absorption is a consequence of co-occurrence between sparse and dense features. If a dense feature like “starts with letter D” always co-occurs with a more sparse feature like “dogs”, the SAE can increase sparsity by absorbing the “starts with D” feature into a “dogs” latent. We explore this further in Appendix A.7, where we show that feature absorption occurs when training an SAE in a toy setting with features that co-occur together.

It remains to be seen if we can predict or identify expected instances where a feature “should have activated” but does not activate due to absorption. One promising direction is meta-SAEs, a novel method for decomposing SAE latents and may decompose absorbed features (Bussmann et al., 2024). One interpretation of our results is that competition may exist between “latents” and “meta-latents” for activation on particular examples and that re-allocation of examples between SAE latents enables SAEs to interpolate between different possible decompositions.

Future Work A primary goal of future work should be to secure further external validity of our findings. This could include finding examples of feature absorption in SAEs trained on other models, with other architectures, or finding examples of feature absorption unrelated to character identification. We expect it should be possible to demonstrate feature absorption in a toy model by mixing dense features with sparse features that always co-occur with these dense features.

We hope as well this investigation may lead to research into solutions, particularly those involving Meta-SAEs (Bussmann et al., 2024), to solve or mitigate feature absorption. Another possible solution may be attribution dictionary learning (Olah et al., 2024a).

Limitations Our feature absorption metric requires having ground-truth knowledge of true labels to first train a LR probe, whereas many features of interest in a LLM lack such clear-cut ground-truth labels. Our metric uses ablation effect to ensure absorbed features causally mediate model behavior, but therefore cannot be easily used in final model layers.

REFERENCES

- 486
487
488 Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier
489 probes, 2017. URL <https://openreview.net/forum?id=ryF7rTqgl>.
- 490 Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick
491 Turner, Cem Anil, Carson Denison, Amanda Askell, et al. Towards monosemanticity: Decom-
492 posing language models with dictionary learning. *Transformer Circuits Thread*, 2, 2023.
- 493 Bart Bussmann, Michael Pearce, Patrick Leask, Joseph Bloom, Lee Sharkey,
494 and Neel Nanda. Showing sae latents are not atomic using meta-saes, 2024.
495 URL [https://www.lesswrong.com/posts/TMAmHh4DdMr4nCSr5/
496 showing-sae-latents-are-not-atomic-using-meta-saes](https://www.lesswrong.com/posts/TMAmHh4DdMr4nCSr5/showing-sae-latents-are-not-atomic-using-meta-saes).
- 497
498 Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, Robert Huben, and Lee Sharkey. Sparse
499 autoencoders find highly interpretable features in language models. In *The Twelfth International
500 Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?
501 id=F76bwRSLeK](https://openreview.net/forum?id=F76bwRSLeK).
- 502 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
503 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony
504 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
505 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,
506 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
507 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
508 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
509 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
510 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
511 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
512 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah
513 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan
514 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
515 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy
516 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
517 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
518 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,
519 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der
520 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,
521 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-
522 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,
523 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
524 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur
525 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-
526 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
527 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
528 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-
529 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
530 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,
531 Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
532 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney
533 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,
534 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,
535 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-
536 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,
537 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,
538 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre
539 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha
Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay
Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda
Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew
Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita

- 540 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh
541 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De
542 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-
543 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina
544 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,
545 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,
546 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
547 Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,
548 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
549 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco
550 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella
551 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory
552 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,
553 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-
554 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,
555 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer
556 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe
557 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie
558 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun
559 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal
560 Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,
561 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian
562 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,
563 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-
564 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel
565 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-
566 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-
567 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,
568 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,
569 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,
570 Piotr Dollár, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,
571 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,
572 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,
573 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-
574 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-
575 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang
576 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
577 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,
578 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,
579 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-
580 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,
581 Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu
582 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-
583 stable, Xiaocheng Tang, Xiaofang Wang, Xiaoqian Wu, Xiaolan Wang, Xide Xia, Xilun Wu,
584 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
585 Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuze He, Zach Rait, Zachary DeVito, Zef
586 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024.
587 URL <https://arxiv.org/abs/2407.21783>.
- 585 Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. Amnesic probing: Behavioral
586 explanation with amnesic counterfactuals. *Transactions of the Association for Computational*
587 *Linguistics*, 9:160–175, 2021.
- 588 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,
589 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposi-
590 tion. *arXiv preprint arXiv:2209.10652*, 2022.
- 591 Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual asso-
592 ciations in auto-regressive language models. In *Proceedings of the 2023 Conference on Empirical*
593 *Methods in Natural Language Processing*, pp. 12216–12235, 2023.

- 594 Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsi-
595 mas. Finding neurons in a haystack: Case studies with sparse probing. *Transactions on Machine*
596 *Learning Research*, 2023. ISSN 2835-8856. URL [https://openreview.net/forum?](https://openreview.net/forum?id=JYs1R9IMJr)
597 [id=JYs1R9IMJr](https://openreview.net/forum?id=JYs1R9IMJr).
598
- 599 Curt Tigges Joseph Bloom and David Chanin. Saelens. [https://github.com/jbloomAus/](https://github.com/jbloomAus/SAELens)
600 [SAELens](https://github.com/jbloomAus/SAELens), 2024.
- 601 Adam Karvonen, Benjamin Wright, Can Rager, Rico Angell, Jannik Brinkmann, Logan Riggs
602 Smith, Claudio Mayrink Verdun, David Bau, and Samuel Marks. Measuring progress in diction-
603 ary learning for language model interpretability with board game models. In *ICML 2024*
604 *Workshop on Mechanistic Interpretability*, 2024.
- 605 Ayush Kaushal and Kyle Mahowald. What do tokens know about their characters and how do they
606 know it? *arXiv preprint arXiv:2206.02608*, 2022.
607
- 608 Connor Kissane, Robert Krzyzanowski, Joseph Isaac Bloom, Arthur Conmy, and Neel Nanda. Inter-
609 preting attention layer outputs with sparse autoencoders, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2406.17759)
610 [abs/2406.17759](https://arxiv.org/abs/2406.17759).
- 611 Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Watten-
612 berg. Emergent world representations: Exploring a sequence model trained on a synthetic task.
613 *ICLR*, 2023.
614
- 615 Tom Lieberum, Senthoran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant
616 Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma Scope: Open Sparse
617 Autoencoders Everywhere All At Once on Gemma 2, August 2024.
- 618 Johnny Lin and Joseph Bloom. Analyzing neural networks with dictionary learning, 2023. URL
619 <https://www.neuronpedia.org>. Software available from neuronpedia.org.
620
- 621 Samuel Marks, Can Rager, Eric J. Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller.
622 Sparse feature circuits: Discovering and editing interpretable causal graphs in language mod-
623 els. *Computing Research Repository*, arXiv:2403.19647, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2403.19647)
624 [abs/2403.19647](https://arxiv.org/abs/2403.19647).
- 625 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associ-
626 ations in GPT. *Advances in Neural Information Processing Systems*, 36, 2022. arXiv:2202.05262.
627
- 628 Neel Nanda. Actually, othello-gpt has a linear emergent world model, mar 2023. URL;
629 <https://neelnanda.io/mechanistic-interpretability/othello>, 2023.
- 630 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.
631 Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
632
- 633 Chris Olah, Adly Templeton, Trenton Bricken, and Adam Jermyn. April update. [https:](https://transformer-circuits.pub/2024/april-update/index.html)
634 [//transformer-circuits.pub/2024/april-update/index.html](https://transformer-circuits.pub/2024/april-update/index.html), 2024a. URL
635 <https://transformer-circuits.pub/2024/april-update/index.html>.
- 636 Chris Olah, Nicholas Turner, Adam Jermyn, and Joshua Batson. July update. [https:](https://transformer-circuits.pub/2024/july-update/index.html)
637 [//transformer-circuits.pub/2024/july-update/index.html](https://transformer-circuits.pub/2024/july-update/index.html), 2024b. URL
638 <https://transformer-circuits.pub/2024/july-update/index.html>.
639
- 640 Bruno A Olshausen and David J Field. Sparse coding with an overcomplete basis set: A strategy
641 employed by v1? *Vision research*, 37(23):3311–3325, 1997.
- 642 Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry
643 of large language models. In *Forty-first International Conference on Machine Learning*, 2024.
644 URL <https://openreview.net/forum?id=UGpGkLzwpP>.
645
- 646 Senthoran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János
647 Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse
autoencoders. *arXiv preprint arXiv:2407.14435*, 2024.

648 Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and
649 use interpretable models instead. *Nature Machine Intelligence*, pp. 206–215, 2019.

650
651 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In
652 *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.

653 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen,
654 Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L
655 Turner, Callum McDougall, Monte MacDiarmid, Alex Tamkin, Esin Durmus, Tristan Hume,
656 Francesco Mosconi, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson,
657 Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Ex-
658 tracting interpretable features from claude 3 sonnet. [https://transformer-circuits.
659 pub/2024/scaling-monosemanticity/](https://transformer-circuits.pub/2024/scaling-monosemanticity/), May 2024. Accessed on May 21, 2024.

660 Matthew Watkins and Joseph Bloom. Linear encoding of character-
661 level information in gpt-j token embeddings, 2023. URL
662 [https://www.lesswrong.com/posts/GyaDCzsyQgc48j8t3/
663 linear-encoding-of-character-level-information-in-gpt-j](https://www.lesswrong.com/posts/GyaDCzsyQgc48j8t3/linear-encoding-of-character-level-information-in-gpt-j).

664
665 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
666 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,
667 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai,
668 Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng
669 Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai
670 Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan
671 Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang
672 Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2
673 technical report. *arXiv preprint arXiv:2407.10671*, 2024.

674 675 A APPENDIX

676 677 A.1 GLOSSARY OF TERMS

678
679 **Sparse Autoencoders (SAEs):** Neural networks trained to reconstruct their input while enforcing
680 sparsity in their hidden layer. In the context of this paper, SAEs are used to decompose the dense
681 activations of language models into more interpretable features.

682
683 **SAE error term:** When inserting a SAE into the computation path of the model, errors in SAE
684 reconstruction will propagate to later parts of the model and can change the model output. We refer
685 to the error as the SAE error term, and corresponds to the difference between the SAE output and
686 the original SAE input activation. Marks et al. (2024) introduced the idea of adding this error term
687 back to the SAE output to ensure that the SAE does not change model output.

688
689 **Latent:** We refer to neurons in the hidden layer of a SAE as latents to avoid overloading the term
690 “feature”. This is in contrast to earlier work which used the term “feature” to refer to both human-
691 interpretable concepts and SAE hidden layer neurons.

692
693 **Feature:** We use the term “feature” to refer to an idealized human-interpretable concept that the
694 model represent in its activations and which a SAE latent may or may not represent.

695
696 **Monosemantic:** Referring to a feature or representation that corresponds to a single, clear seman-
697 tic concept. In the context of SAEs, a monosemantic feature would ideally capture one interpretable
698 aspect of the input.

699
700 **Interpretable:** A latent being interpretable is not well defined in the field, making it difficult to
701 ensure that different authors mean the same thing when referring to SAE interpretability. When
we refer to a SAE latent as being interpretable in this work, we mean that it should behave in line
with how it appears to behave after inspecting its activation patterns. If an SAE latent appears

702 to track a feature X by a reasonable inspection of its activations but has subtle deviations from this
703 behavior in reality, we say this is not interpretable. We thus measure interpretability via classification
704 performance when a latent appears to be a classifier over some feature.
705

706 **Feature dashboard:** A dashboard showing activation patterns and max-activating examples for a
707 SAE latent. Feature dashboards are commonly used to interpret the behavior of an SAE latent.
708

709 **Neuronpedia:** A platform, <https://neuronpedia.org>, which hosts feature dashboards for
710 popular SAEs (Lin & Bloom, 2023).
711

712 **Token-aligned latent:** A latent which seems to roughly fire on variants of the same token. For
713 instance, a “Snake” token-aligned latent may fire on the tokens “Snake”, “SNAKE”, “_snakes”,
714 etc...
715

716 **Feature splitting:** A phenomenon in SAEs introduced by Bricken et al. (2023), where a SAE
717 latent tracks a general feature in a narrow SAE, but splits into multiple more specific SAE latents in
718 a wider SAE. For instance, a latent tracking “starts with L” in a narrow SAE may split into a latent
719 tracking “starts with capital L” and a latent tracking “starts with lowercase L” in a wider SAE.
720

721 **Feature absorption:** A problematic form of feature splitting where a SAE latent appears to track
722 an interpretable feature, but that latent has seemingly arbitrary exception cases where it fails to fire.
723 Instead, an approximately token-aligned feature “absorbs” the feature direction and fires in place of
724 the main latent.
725

726 **Circuit:** In the context of neural network interpretability, a circuit refers to a subgraph of neurons
727 or features within a neural network that work together to perform a specific function or computation.
728 The study of circuits aims to understand how different components of a neural network interact to
729 process information and produce outputs.

730 **Linear probe:** A simple linear classifier (typically logistic regression) trained on the hidden ac-
731 tivations of a neural network to predict some property or task. Used to assess what information is
732 linearly decodable from the network’s representations.
733

734 **K-sparse probing:** A variant of linear probing where only the k most important features (as de-
735 termined by some selection method) are used to train the probe. This helps identify which specific
736 neurons or features are most relevant for a given task.
737

738 **Ablation study:** An experimental method where a component of a system (in this case, a neuron or
739 feature in a neural network) is removed or altered to observe its effect on the system’s performance.
740 This helps determine the causal importance of the component.
741

742 **Integrated gradients (IG):** An attribution method that assigns importance scores to input features
743 by accumulating gradients along a path from a baseline input to the actual input. In this paper, it’s
744 used as an approximation technique for ablation studies.
745

746 **In-context learning (ICL):** A paradigm where a language model is given examples of a task
747 within its input prompt, allowing it to adapt to new tasks without fine-tuning. Often used with
748 few-shot learning techniques.
749

750 **Residual stream:** In the context of transformer architectures, the residual stream refers to the
751 main information flow that bypasses the self-attention and feed-forward layers through residual
752 connections.
753

754 **Logits:** The raw, unnormalized outputs of a neural network’s final layer, before any activation
755 function (like softmax) is applied. In language models, logits typically represent the model’s scores
for each token in the vocabulary.

Activation patching: An interpretability technique where activations at specific locations in a neural network are replaced or modified to observe the effect on the network’s output. This helps in understanding the causal role of different parts of the network in producing its final output.

A.2 HOW GOOD IS GEMMA-2 ON CHARACTER IDENTIFICATION TASKS?

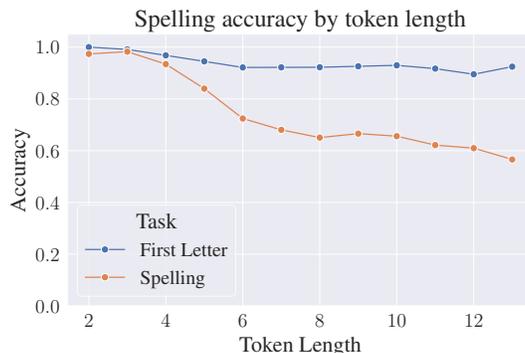


Figure 8: Baseline performance for Gemma-2-2B on first-letter identification and full-token spelling by token length.

We evaluate how well can Gemma-2-2B identify the first letter or all the letters in a token (spelling the full token). We evaluate the accuracy of the model on all tokens in the LR probe validation set with a prompt containing 10 in-context examples selected at random from the full vocabulary. Our results are shown in Figure 8.

We see that performance on the first-letter identification task is high throughout token length, while the full-word spelling performance decreases as the length of the token increases.

A.3 INTERVENING ON THE FIRST LETTER

If the model is using the identified SAE latents for predicting the first letter we should also be able to change what first letter it predicts just by changing the activations. For this experiment we use the SAE latents most cosine similar with the LR probe for the true first letter and for a new randomly selected letter. We take the intermediate activations of Gemma-2-2B in the residual stream and encode them using the SAE. Then we zero out the activation of the SAE latent associated with the original letter and change the activation of the SAE latent associated with the new letter into the average activation it has on tokens starting with this new letter.

Editing works better with latents from the narrower 16k SAE compared to the 65k, with the best L0s in the 75-150 range. This corresponds to the observed pattern of these SAE latents having higher F1 scores for classification. We report the results in Figure 9. The best SAEs on the layers 7-9 can achieve a substantial replacement, but note that the averages hide variance across individual tokens, where some get edited completely and others get unaffected. The edit success also varies based on the true first letter and the random new letter; for illustration we show a breakdown by letter for two specific SAEs in layer 7 in Figure 10.

A.4 PROBE COSINE SIMILARITY VS K=1 SPARSE PROBING

The first step when searching for a SAE feature that acts as a first-letter classifier involves searching for SAE feature which best acts as a classifier. In Figure 2, we achieve this by first training a LR probe on the first-letter task and using cosine similarity between that probe and the SAE encoder to find the best feature for the first-letter task. We also investigated using k-sparse probing with k=1 to select the best SAE feature instead. This involves training a linear probe with L1 loss and selecting the feature with the highest positive weight from the probe.

We find that both k=1 sparse probing yield nearly identical results, as seen in Figures 11 and 12. Additionally Figure 13 shows the cosine similarity of the LR probe with each SAE feature by letter



(a) Comparing success in editing out the true first letter and making the model predict a randomly selected new letter across layers 0-9 for all 16k and 65k Gemma Scope SAEs.

(b) Comparing the edit success with the top SAE feature across all L0s for 16k and 65k widths across layers 0-9. The best performance seems to be occurring for L0 between 75 and 150.

Figure 9: Comparison of Edit success by Layer and L0

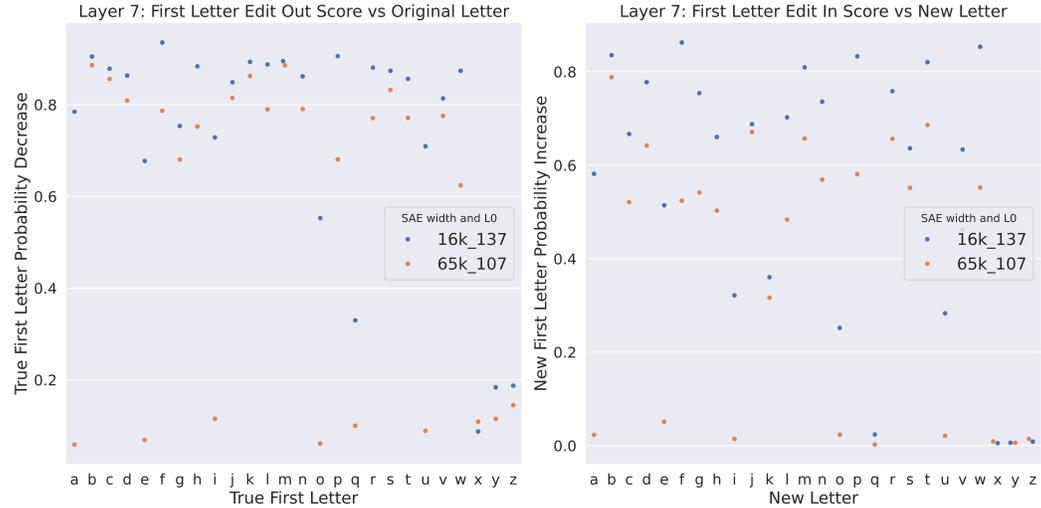


Figure 10: Comparing the edit success broken down by the letter at layer 7 for two SAEs; SAE width 16,000 and L0 of 137 and SAE width 65,000 and L0 of 107. For each original letter we draw a sample of 100 tokens and average the decrease in probability of the correct first letter and increase in probability of a new random letter.

for the canonical Gemma Scope layer 0 16k width SAE. In most cases there is an obvious probe-aligned feature. Likely any reasonable method of feature selection will find the same feature for these cases. We thus decided to use cosine similarity between the SAE encoder and a LR probe as our selection criteria for single SAE features as this is a simpler metric and less computationally intensive to compute.

A.5 PRECISION, RECALL, AND F1 SCORE FOR THE FIRST-LETTER TASK

We evaluated precision, recall, and F1 score for the first-letter classification task, and found that the precision and recall vary depending on the L0 of the SAE. Low L0 SAEs learn high precision, low recall features, while high L0 SAEs learn low precision, high recall features. These results are shown in Figure 14. We thus chose to use F1 score as our core metric in this paper to balance precision and recall as many of the SAEs we tested have extreme values in either precision or recall.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

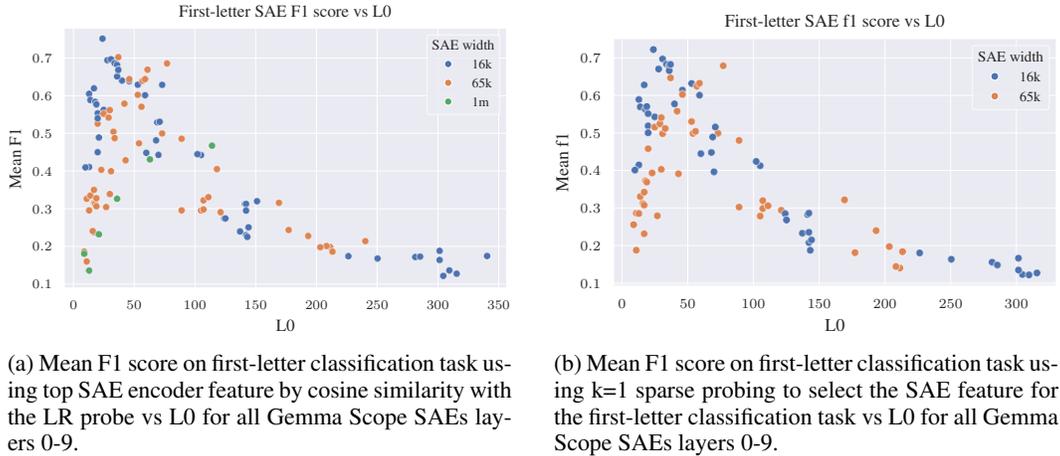


Figure 11: Comparison of LR probe cosine similarity and k=1 sparse probing vs l0

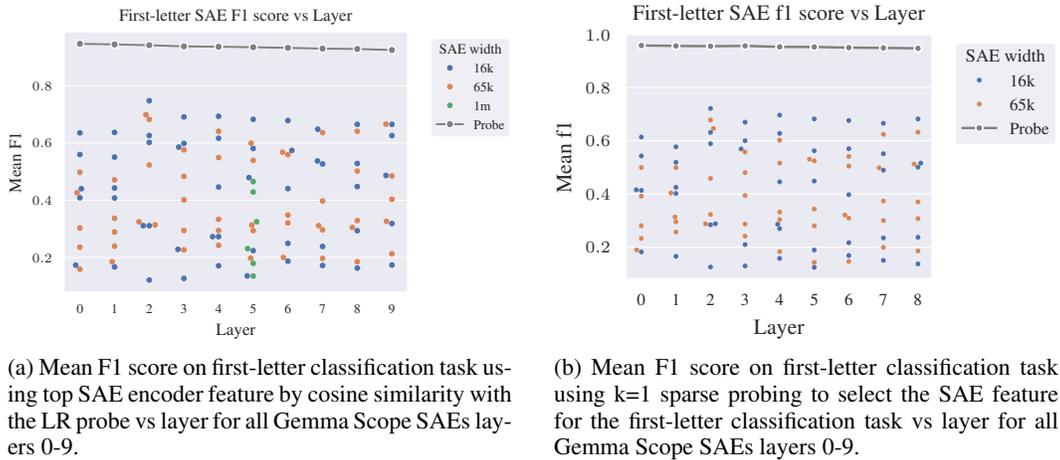
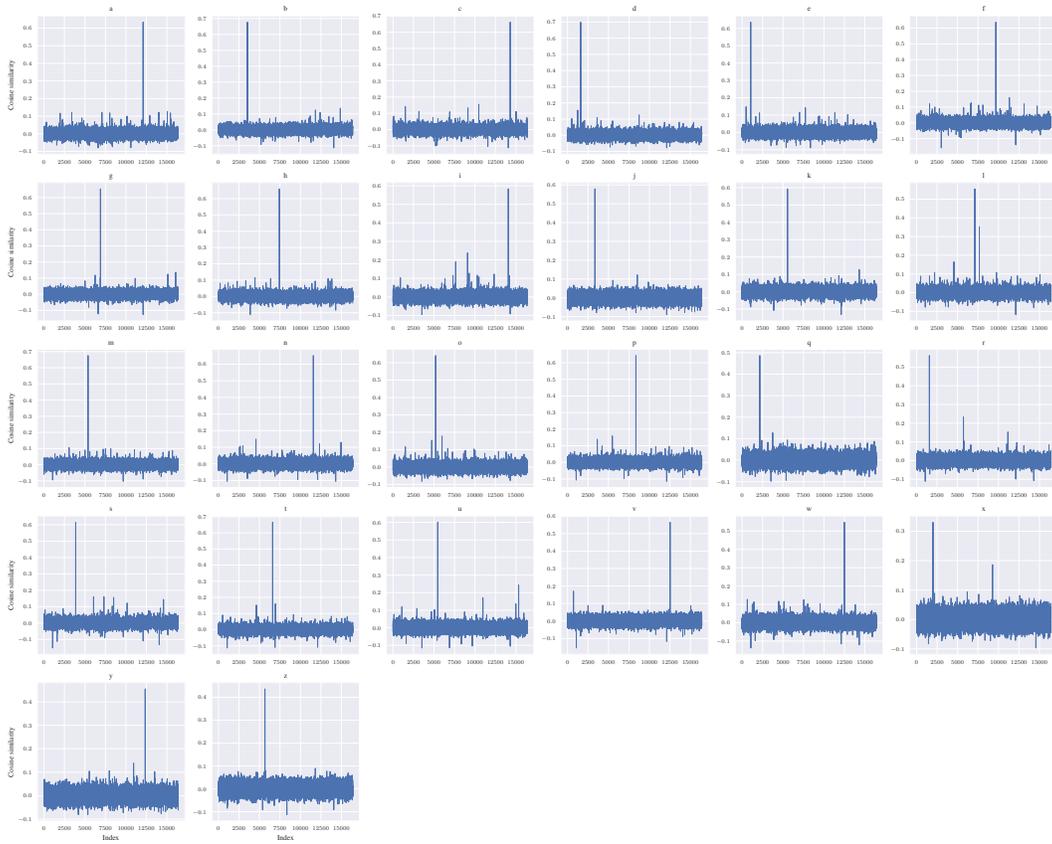


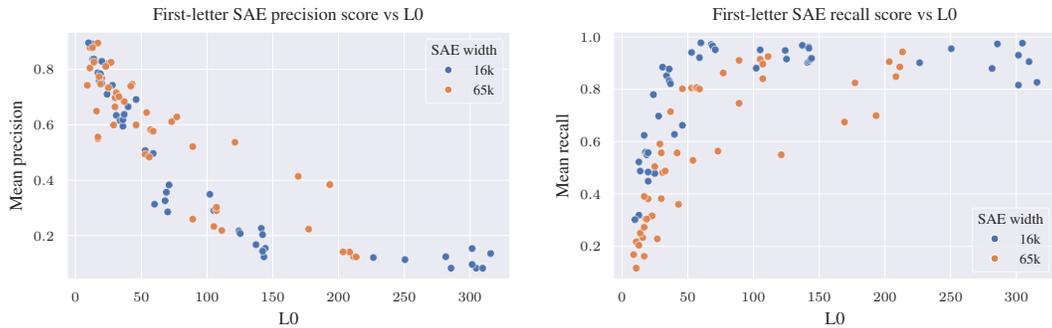
Figure 12: Comparison of LR probe cosine similarity and k=1 sparse probing vs layer

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944



945 Figure 13: Decoder cosine similarities with the LR probe by letter, Gemma Scope 16k layer 0
946 $l_0=105$. Most letters have one or two obvious SAE features which align with the probe.
947

948
949
950
951
952
953
954
955
956
957
958
959



960 (a) Mean precision on first-letter classification task vs L_0 for all Gemma Scope SAEs layers 0-9. Features
961 are selected via $k=1$ sparse probing
962 (b) Mean recall on first-letter classification task vs L_0
963 for all Gemma Scope SAEs layers 0-9. Features are
964 selected via $k=1$ sparse probing

965 Figure 14: Comparison of precision and recall vs l_0

966
967
968
969
970
971

While it may appear that there is an optimal L_0 from looking at aggregate statistics across letter, we find that breaking down the F1 vs L_0 plot by letter reveals that the optimal L_0 appears different for different letters, with low frequency letters like z actually having the best F1 score at the lowest L_0 , while other letters instead have an optimal L_0 around 30-50. Figure 15 shows these results broken down by letter.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999

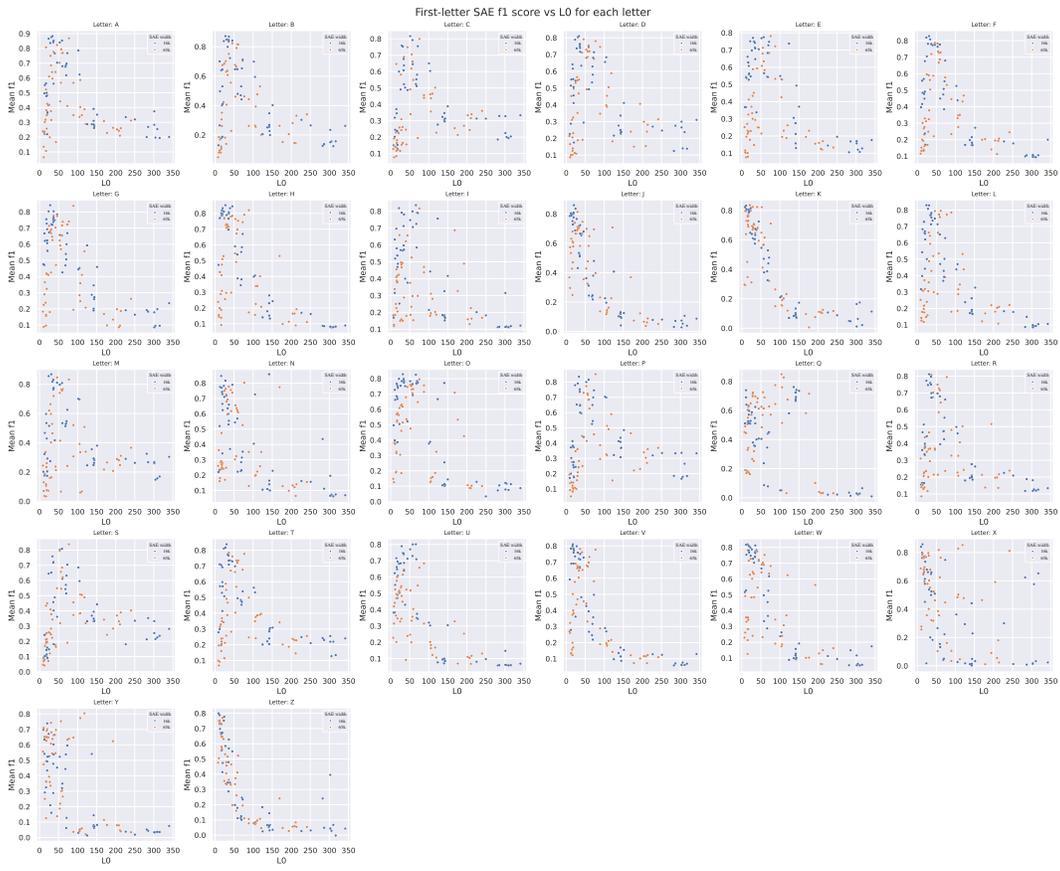


Figure 15: F1 vs LO by letter. SAE features are picked using k=1 sparse probing.

1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010

A.6 SAE TRAINING

We train SAEs on the first 8 layers of Qwen2 0.5B (Yang et al., 2024) and Llama 3.2 1B (Dubey et al., 2024) using the SAELens library (Joseph Bloom & Chanin, 2024). The SAEs are all trained with identical hyperparameters of L1 coefficient of 2.5 and 500M tokens. The Qwen2 0.5B SAEs all have LO between 25 and 50 and explained variance between 0.77 and 0.83. The Llama 3.2 1B SAEs have LO between 27 and 110, and explained variance between 0.74 and 0.89.

A.7 TOY MODELS OF FEATURE ABSORPTION

We hypothesize that absorption is due to feature co-occurrence combined with the SAE maximizing sparsity. When two features co-occur, for instance "starts with S" and "short", the SAE can increase sparsity by merging the "starts with S" feature direction into a latent tracking "short" and then simply not fire the main "start with S" latent. This means firing one feature instead of two, and thus increasing sparsity.

We show that feature co-occurrence does indeed cause absorption by constructing a toy setting with four true features and an SAE with four latents.

1020
1021
1022
1023
1024
1025

Setup Our initial setup consists of 4 true features, each randomly initialized into orthogonal directions with a 50 dimensional representation vector and unit norm. We control the base firing rates of each of the 4 true features. Unless otherwise specified, the feature fires with magnitude 1.0 and stdev 0.0. We train a SAE with 4 latents to match the 4 true features using SAELens (Joseph Bloom & Chanin, 2024). The SAE uses L1 loss with l1 coefficient 3e-5, and learning rate 3e-4. We train on 100,000,000 activations. Our 4 true features have the firing rates shown in Table 2.

	Feature 0	Feature 1	Feature 2	Feature 3
Firing rate	0.25	0.05	0.05	0.05

Table 2: Base firing rate for features in toy experiment.

We use this setup for the following reasons:

- This is a very easy task for a SAE, and it should be able to reconstruct these features nearly perfectly.
- Using fully orthogonal features lets us see exactly what the L1 loss term incentivizes without worrying about interference from superposition.

Independently firing features When the true features fire independently, we find that the SAE is able to perfectly recover these features as shown in Figure 16.



Figure 16: When features fire independently, the SAE learns exactly one latent per feature, and the decoder perfectly reconstructs the feature.

We see the cosine similarity between the true features and the learned encoder, and likewise with the true features and the decoder. The SAE learns one latent per true feature. The decoder representations perfectly match the true feature representations, and the encoder learns to perfectly segment out each feature from the other features.

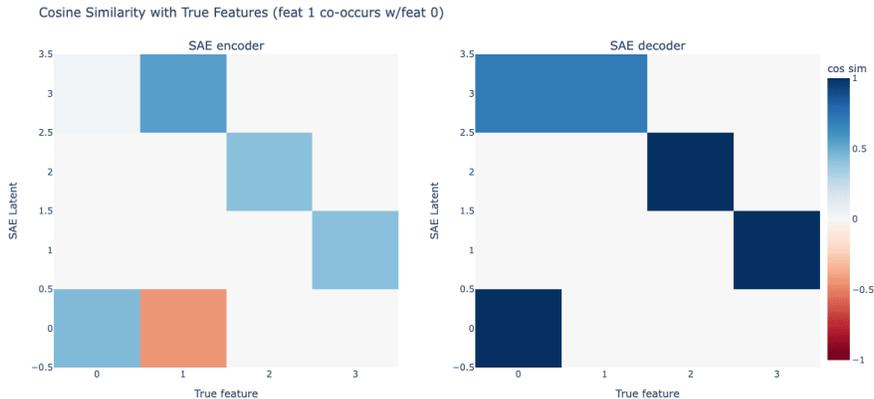
Co-occurrence causes absorption Next, we modify the firing pattern of feature 1 so it fires only if feature 0 also fires. However, we keep the overall firing rate of feature 1 the same as before, firing in 5% of activations. Features 2 and 3 remain independent.

Figure 17 shows the encoder and decoder cosine similarities with the true features in the co-occurrence setup. Here, we see a clear example of feature absorption. Latent 0 has learned a perfect representation of feature 0, but the encoder has a hole in its recall. Latent 0 fires if feature 0 is active but not feature 1. This is exactly the sort of gerrymandered feature firing pattern we saw in real SAEs for the starting letter task - the encoder has learned to stop the latent firing on specific cases where it looks like it should be firing. In addition, we see that latent 3, which tracks feature 1, has absorbed the feature 0 direction. This results in latent 3 representing a combination of feature 0 and feature 1. We see that the independently firing features 2 and 3 are untouched - the SAE still learns perfect representations of these features.

A.8 METRIC CHOICE FOR ABLATION STUDIES

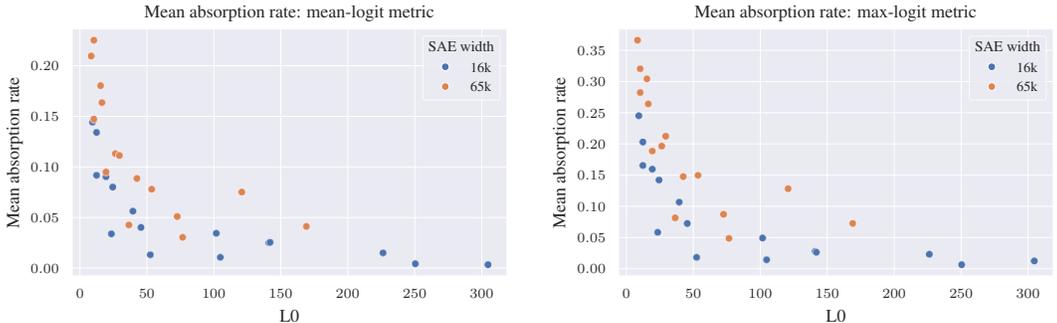
To determine the causal effect of SAE latents on the first-letter identification task, we use a metric, m , which measures the logit of the correct letter minus the mean logit of all incorrect letters. Our

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094



1095 Figure 17: When features 0 and 1 co-occur, we see absorption in the SAE encoder and decoder
1096 latents which track features 0 and 1.

1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108



1109 (a) Absorption rate using the mean version of the absorption metric, Gemma Scope layers 0-3.
1110 (b) Absorption rate using the max version of the absorption metric, Gemma Scope layers 0-3.

1111 Figure 18: Comparison of absorption rates using the max and mean versions of the absorption
1112 metric.

1113
1114

1115 metric is defined below, where g refers to the final token logits, L is the set of uppercase letters, and
1116 y is the uppercase letter that is the correct starting letter:

1117
1118
1119
1120

$$m = g[y] - \frac{1}{|L| - 1} \sum_{l \in \{L \setminus y\}} g[l]$$

1121
1122
1123
1124

This metric is chosen to detect changes in the confidence of the model in predicting the correct letter
relative to the mean reference class of other letters. This should capture changes in the model’s
confidence in predicting the correct logit.

1125
1126
1127

This is not the only metric that could be chosen, and an argument can be made that we should
subtract the max of all incorrect letter logits rather than the mean of all incorrect letter logits. The
max form of this version of the metric is shown below:

1128
1129
1130
1131

$$m_{max} = g[y] - \max_{l \in \{L \setminus y\}} g[l]$$

1132
1133

This second form using a max can also account for the case where the logits of the model shift from
being confident in the correct answer to instead being confident in an incorrect answer while leaving
the logits of the correct answer the same.

In practice, we expect that ablating an absorbing latent should cause the model to become less confident in the correct answer, so the difference between these two forms of the metric should yield similar results.

We calculate the mean absorption rate for Gemma Scope SAEs layers 0-3 in Figure 18 using both versions of this metric. The overall shape of the curve is nearly identical between these two choices of metrics. The mean version of the metric, which is used in this paper, results in a slightly more conservative estimate of absorption rate.

We consider our absorption score to be a rough estimate of the true absorption rate and thus consider either the mean or the max version of the logit diff metric to be valid for evaluating absorption.

A.9 CAUSAL INTERVENTIONS AND ABSORPTION

In this work, we rely on causal interventions like ablation experiments to verify that SAE latents have a causal impact on model behavior. In these experiments for spelling tasks, we set up an ICL prompt to elicit spelling information from the model, for instance the ICL prompt below:

```
tartan has the first letter: T
mirth has the first letter: M
dog has the first letter:
```

In this ICL prompt, we would apply an SAE and train LR probes on the `_dog` token position, and expect that the model will output the token `_D`. When we intervene on the `_dog`, we can track the causal changes to model outputs by applying a metric to the output logits, e.g. checking how our intervention increases or decreases the `_D` logit relative to other letters.

We use these interventions as part of our absorption metric to ensure that when we claim that “absorption” is occurring, we verify that the absorbing feature has a causal impact on model outputs. This is stronger evidence than only noting a cosine similarity between the absorbing feature, but this means that our absorption metric cannot classify absorption at later model layers.

During a LLM forward pass, the model first collects relevant information on a token in that token position, and attention heads then move relevant information from earlier tokens to later tokens (Geva et al., 2023; Meng et al., 2022). If we assess ablation effect at layers after which model attention has already pulled relevant information from the subject tokens into the final output token, the ablation effect will be 0. For Gemma 2 2B on the first-letter spelling task, we find this movement of first-letter spelling information occurs around layer 18.

Figure 19 shows an activation patching experiment (Meng et al., 2022) on a sample first-letter spelling prompt. In this experiment, we see that near layer 18 the model moves first-letter spelling information from the subject token to the prediction token.

As a result, our feature absorption metric will not function past layer 18 in Gemma 2 2B, and we thus focus on layers 0-17 for our analysis of feature absorption. We believe that feature absorption is still occurring in SAEs past layer 18, but we lose the ability to make causal claims that the absorbing features are used by the model to make predictions. Given that this paper is trying to highlight the existence of feature absorption, we felt it is more important to have a metric which is robust and has the backing of causal analysis but which cannot be used at all model layers. Future work may make a different trade-off and choose a feature absorption metric which can work at all model layers, for instance relying only on cosine similarity between absorbing features and a LR probe to determine absorption.

A.10 ADDITIONAL PLOTS

In this section, we include additional plots that are too large to fit in the main body of the paper.

A.11 FEATURE DASHBOARDS

We include feature dashboard screenshots from Neuronpedia for some prominent latents mentioned in this work. Figure 21 shows a dashboard for Gemmascope layer 3, latent 1085, which is a token-aligned latent firing on variations of the word `_short` and we find absorbs the “starts with S”

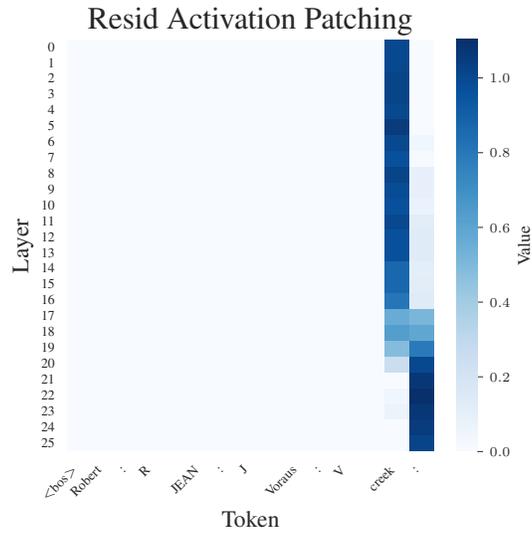


Figure 19: Residual stream attribution patching for a sample first-letter spelling prompt, Gemma 2 2B. After around layer 18, model attention moves the relevant spelling information from the source token to the prediction location.

direction. Figure 22 shows latent 6510 from the same layer which should be the main “starts with S” latent.

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

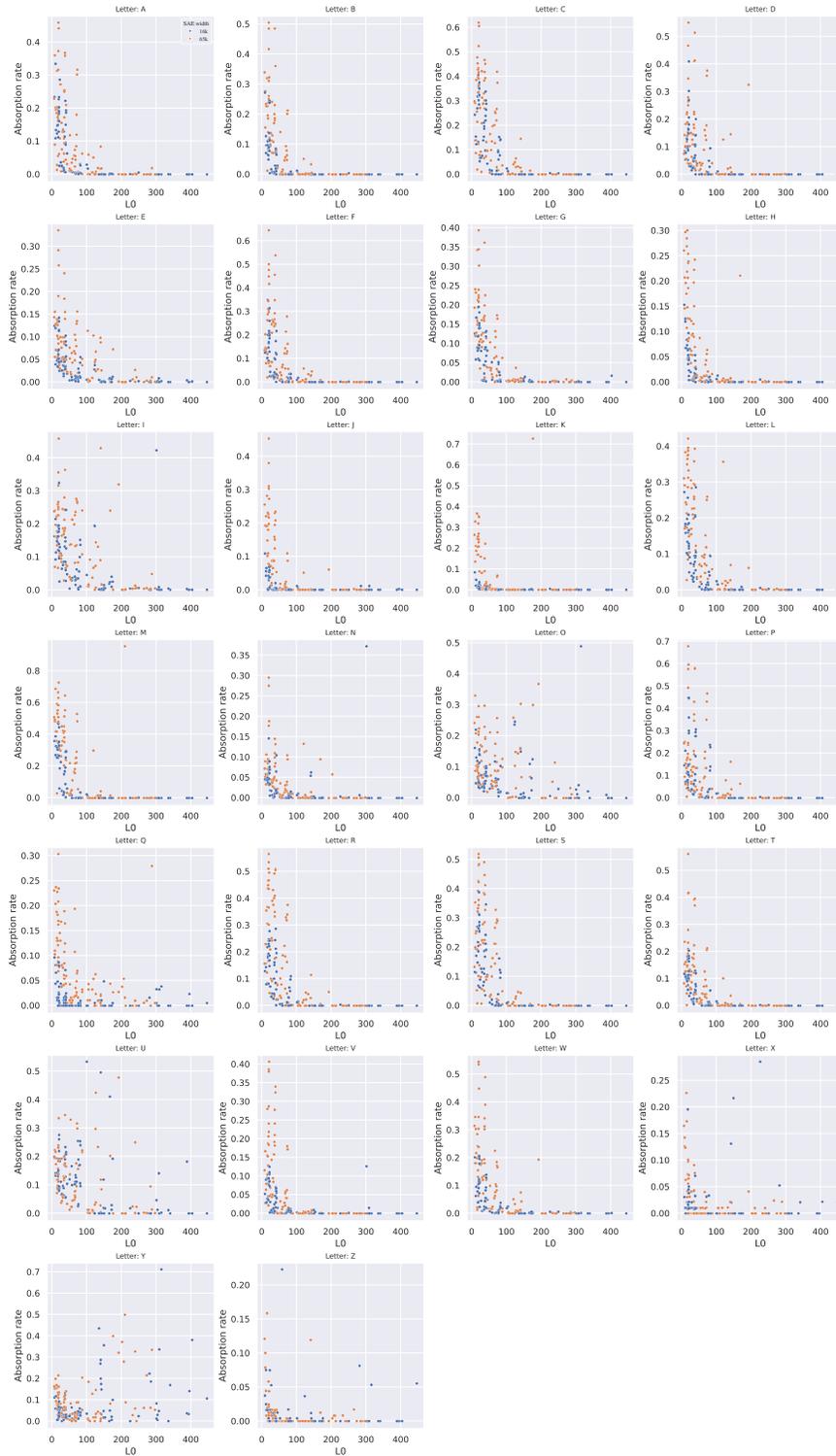


Figure 20: Absorption rate vs L0 by letter, layers 0-17. We see a wide variance in which letters are absorbed by which SAEs.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317

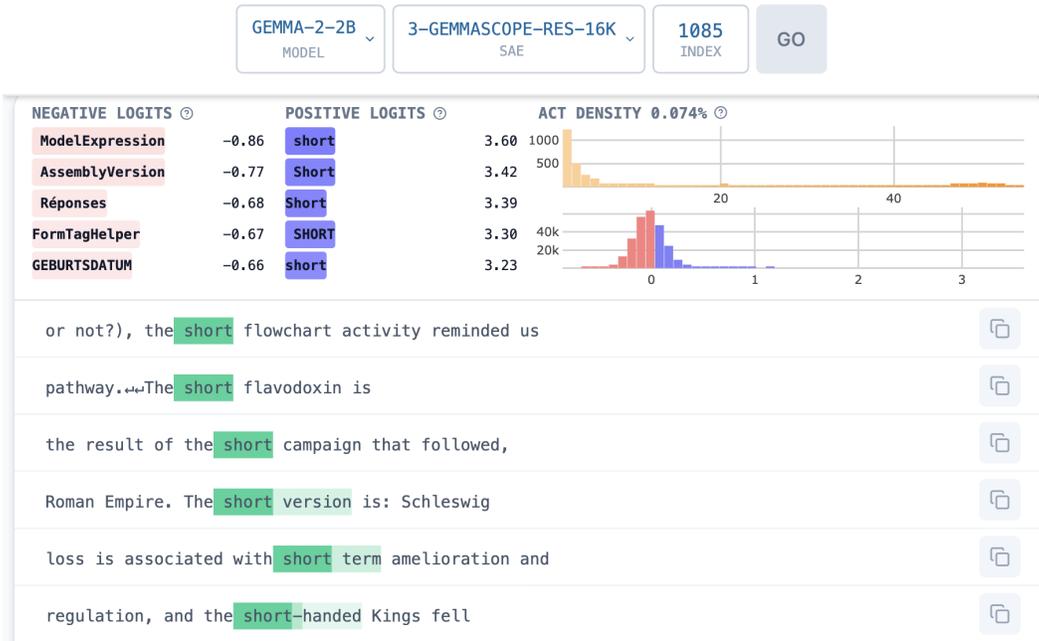


Figure 21: Neuronpedia dashboard for Gemma Scope layer 3, latent 1085. This latent is a token-aligned latent for `_short` tokens. This latent absorbs the “starts with S” direction.

1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

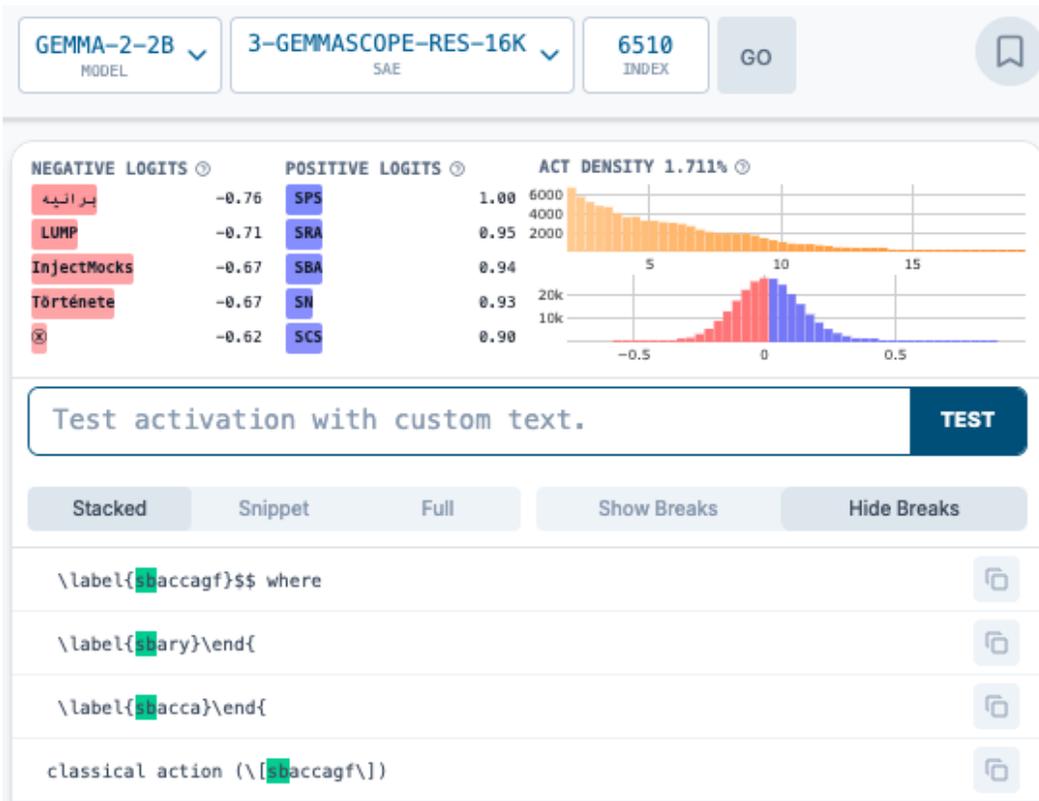


Figure 22: Neuronpedia dashboard for Gemma Scope layer 3, latent 6510. This latent should be the main “starts with S” latent.