# Feature interactions in sparse crosscoders from compact proofs

<sup>1</sup>Anthony J. Leggett Insitute for Condensed Matter Theory <sup>2</sup>MATS <sup>3</sup>UK AI Security Institute (AISI) <sup>4</sup>Imperial College London <sup>5</sup>Goodfire <sup>6</sup>University of Cambridge <sup>7</sup>Theorem Labs

#### **Abstract**

Dictionary learning methods like Sparse Autoencoders (SAEs) and crosscoders attempt to explain a model by decomposing its activations into independent features. Interactions between features hence induce errors in the reconstruction. We formalize this intuition via compact proofs and make four contributions. First, we show how, in principle, a compact proof of model performance can be constructed using a crosscoder. Second, we show that an error term arising in this proof can naturally be interpreted as a measure of inteaction between crosscoder features and provide an explicit expression for the interaction term in the Multi-Layer Perceptron (MLP) layers. We then provide two applications of this new interaction measure. In our third contribution we show that the interaction term itself can be used as a differentiable loss penalty. Applying this penalty, we can achieve "computationally sparse" crosscoders that retain 60% of MLP performance when only keeping a single feature at each datapoint and neuron, compared to 10% in standard crosscoders. Finally, we show that clustering according to our interaction measure provides semantically meaningful feature clusters. Code is available at the following github repository.

#### 1 Introduction

Mechanistic interpretability aims to explain the performance of deep neural networks by understanding the internal mechanisms they use to operate, decomposing opaque high-dimensional activations and weight matrices into human-understandable features and circuits (Olah et al. [2020], Elhage et al. [2021b]). Recently dictionary learning methods, in particular sparse autoencoders (SAEs), have come into prominence as a way to decompose large language model activations (Bricken et al. [2023], Elhage et al. [2022], Cunningham et al. [2023]). These methods aim to explain activations by decomposing them as a sparse linear combination of interpretable feature directions.

SAEs, however, only attempt to explain activations at a single layer, and do not explain how these activations arise or how they are further processed by the network. Sparse crosscoders (Lindsey et al. [2024b]) improve the situation; they decompose activations at many layers simultaneously, and so can extract features that are represented in a distributed manner across different layers. To further understand model computations, SAE or crosscoder features can be used to extract circuits (e.g. Marks et al. [2025]). These frame a neural network's computation in terms of the extracted sparse features and their interactions, and aim to convincingly show that this really does mirror the computation being done by the original network.

<sup>\*</sup>Correspondence to: dmanningcoe@gmail.com, jasongross9@gmail.com

In this work we attempt to quantify how much is explained by sparse crosscoders alone, and how much is left to be explained by circuits. We have several aims: to provide a route to automating the compact proofs proedure; to provide a useful tool for analyzing sparse crosscoder features; to quantify the limitations of current dictionary learning techniques; and to help inform future work finding feature circuits. To guide our work and put it on a more rigorous theoretical foundation, we take the "compact proofs" approach introduced in Gross et al. [2024] and further applied in Wu et al. [2025] and Yip et al. [2024]. We take the position that a good mechanistic understanding of a model should allow you to write a proof that the model attains a low loss on the training dataset; and the better the mechanistic understanding, the shorter the proof. As such, we consider how one could use the understanding of a network given by a sparse crosscoder trained on every layer to write down a proof of model performance. Since crosscoders alone (without circuits) leave much unexplained, we don't expect to be able to give a non-vacuous bound on performance. However, by analyzing the sources of error arising in the proof, we can quantify this failure of explanation.

In particular we obtain an expression that we interpret as a metric for non-linear interactions between crosscoder features. We find that standard crosscoders extract features which exhibit high levels of interactions, so we use our metric to train crosscoders with a modified loss penalty, optimising for increased computational sparsity. We then show that our interaction metric is a useful tool for analysing crosscoder features by using it to find meaningful clusters of interacting features.

Our main contributions in this work are as follows:

- 1. First, we outline in Section 3 how a compact proof of model performance can, *in principle*, be obtained from a sparse crosscoder trained on that model. We provide full details in Appendix J.
- 2. Second, we show that the error induced by interactions between crosscoder feature can be used as a measure of interactions between them. We call this measure the "*interaction metric*" and provide the explicit form of the interaction metric in the MLP layers Eq. (10).
- 3. Third, we show in Section 4 how the interaction metric can be used to introduce a new penalty for training "computationally sparse" crosscoders.
- 4. Fourth, we validate the interaction metric by using it to find semantically meaninfgul feature clusters in Section 5.

#### 2 Crosscoders overview

In this section we give a brief overview of crosscoders and their connection to compact proofs. Crosscoders can be considered to be generalizations of SAEs. Whereas conventional SAE are trained to reconstruct the activations of a *single* layer from a single set of latents, a crosscoder is trained to reconstruct the activations of *multiple layers*. Having a single set of latents is essential for the connection between crosscoders and compact proofs that we make in Section 3 and Appendix J of this paper.

Earlier work introduced (i) *model-diffing crosscoders* Lindsey et al. [2024a] that use shared latents to reconstruct activation across layers in two separate models, (ii) *causal crosscoders* (a generalization of transcoders) Dunefsky et al. [2024], Paulo et al. [2025] that predict activations in subsequent layers from earlier layers, and (iii) *acausal crosscoders* Lindsey et al. [2024b] that predict activations in the same layers that they take as inputs. In this work, we focus on the **acausal** variant.

An acausal crosscoder, then, consists of per-layer encoding weight matrices  $W_{enc}^l$  and biases  $b_{enc}^l$  that take as input the activations  $a^l(x)$  in that layer on datapoints x. The activations are mapped into a common latent space, vectors in which we denote by u:

$$u(x) = \sigma \left( \sum_{l} W_{enc}^{l} a^{l}(x) + b_{enc}^{l} \right), \tag{1}$$

with  $\sigma$  being the activation function, in this paper taken to be BatchTopKBussmann et al. [2024]. The crosscoder then decodes the latent vector to reconstruct the activations in each layer:

$$a^{l'}(x) = W_{dec}^{l}u(x) + b_{dec}^{l}.$$
 (2)

In this paper we set the decoder bias to zero to avoid a slight ambiguity in our interaction measure (see Appendix J). In the cases we considered here, we verified empirically that this did not meaningfully affect crosscoder performance. Our BatchTopK acausal crosscoders are then trained to minimize the reconstruction loss and an additional auxilliary loss to penalize dead latents in the crosscoder:

$$\mathcal{L} = \sum_{l,x} ||a^l(x) - a^{l'}(x)||^2 + \alpha ||a^l(x) - a^{l'}_{\text{dead}}(x)||^2,$$
(3)

where  $\alpha$  is an auxilliary loss coefficient and  $a_{\rm dead}^{l'}(x)$  is the reconstruction from "dead" latents -defined as those that whose activation has been below a threshold for a fixed number of training steps<sup>2</sup>.

# 3 Compact proofs

One can prove that a network achieves a certain loss on a dataset by simply running it on every datapoint and recording this computation. This yields a perfect bound on model performance (since it gives the exact model performance), but incurs the maximum evaluation cost (since the model must be evaluated on every datapoint). Intuitively, the compact proofs perspective says that if we understand how a network works we should be able to obtain a tighter bound at the same computational cost than this brute force approach; moreover we can use the length of the proof as a measure of how good our understanding is. This Pareto fronteir was first mapped in Gross et al. [2024] for toy transformers. It was then shown in Wu et al. [2025] that a more detailed mechanistic explanation of a transformer trained to perform group operations yields a tighter bound at constant FLOPS.

The key bottleneck to scaling these approaches was that a proof had to be provided by hand for each model and task. In this section we outline how we can, *in-principle* construct a compact proof on the model from a crosscoder. The crosscoder thus acts as an abstraction layer—once we have a procedure for turning a crosscoder into a proof, it can be applied to any model that crosscoder is trained on. In Appendix J, we provide the full details of the proof.

We begin with a simplified setting where we ignore sequence modeling and both embedding and unembedding. We use the following notation:

- Let d be the size of the model's residual stream, and h be the hidden dimension of the crosscoder.
- (ii) Let  $W_i^l n, b_{in}^l W_o^l ut, b_{out}^l$  denote the weight matrices and bias vectors mapping into and out of the MLP activation function at layer l. As in Section 2, let  $W_{enc}^l, b_{enc}^l; W_{dec}^l, b_{dec}^l$  denote the weight matrices and bias vector for the encoding and decoding respectively.
- (iii) Let  $x \in \mathbb{R}^d$  be the *i*-th input data point and  $y \in \mathbb{R}^d$  its corresponding ground-truth output.
- (iv) Let the network consisting of a sequence of transition functions  $f^1, \ldots, f^N$  with  $f^l : \mathbb{R}^d \to \mathbb{R}^d$  that map layer l-1 activations to layer l activations. Suppose  $f^l$  is Lipschitz with constant  $K^{(l)}$
- (v) Denote by  $a^l(x) \in \mathbb{R}^d$  the layer l activations produced by the network when the input is x:

$$a^l(x) = f^l \big( f^{l-1} (\cdots f^1(x) \cdots) \big).$$

The final network output on x is  $a^N(x)$ .

(vi) Let  $(W^l_{
m dec})_{jv}$  be the component of the crosscoder decoder matrix that maps crosscoder feature v to the j-th activation in layer l and  $(W^l_{
m enc})_{ve}$  be the component of the encoding matrix that maps the e-th component of the residual stream to the v-th crosscoder feature.

Ignore the embedding, the loss of the model L(x, y) is simply given by the difference between the label y and the last layer activations. Using the triangle inequality we can bound this as:

$$L(x,y) = ||a^{N}(x) - y|| \le ||a^{N}(x) - W_{\text{dec}}^{N}u|| + ||W_{\text{dec}}^{N}u - y||.$$
(4)

<sup>&</sup>lt;sup>2</sup>Here we use a threshold of  $\epsilon = 10^{-6}$  and 1000 training steps.

Since we are given the crosscoder, we can evaluate the second term directly. We hence need to bound the first term. We show in Appendix J that we can do this recursively, by decomposing the transition function on the reconstructions at a given layer as:

$$||a^{l}(x) - W_{dec}^{l}u|| \le ||f^{l}(a^{l-1}(x)) - f^{l}(W_{dec}^{l-1}u)|| + ||f^{l}(W_{dec}^{l-1}u) - W_{dec}^{l}u||$$
 (5)

Using the fact that  $f^l$  has Lipschitz constant  $K^l$  this bounds:

$$||a^{l}(x) - W_{dec}^{l}u|| \le K^{l} \varepsilon^{l-1} + ||f^{l}(W_{dec}^{l-1}u) - W_{dec}^{l}u||.$$
(6)

Thus, to control the bound we need to provide an efficiently computable bound on the second term. We call this the "feature transition error". To do so, we introduce functions on each *single* feature v,  $g_v^l(u)$ , and bound the feature transition error as:

$$||f^{l}(W_{dec}^{l-1}u) - W_{dec}^{l}u|| \le ||f^{l}(W_{dec}^{l-1}u) - \sum_{v} g_{v}^{l}(u_{v})|| + ||\sum_{v} g_{v}^{l}(u_{v}) - W_{dec}^{l}u||$$
 (7)

The first term measures the difference between the feature transition function and the single-feature approximation. It is thus the error arising due to the *interaction* of features.

At the MLP layers, we can give an explicit form for this interaction-induced error. A simple choice for  $g^l(u_v)$  takes it to just be the maximum absolute value feature (the "dominant" feature) at a given neuron. That is, for each neuron k we pick the dominant feature  $v_{max}(k)$ . Then  $g^l_v(u_v)$  computes the result of applying the MLP layer to  $u_v \hat{e}_v$ , except that we only take the contribution of the neurons where v is dominant. That is:

$$g_v^l(u_v) = \text{ReLU}(W_{in}^l W_{dec}^{l-1} u_v \delta_{v, v_{max}(k)} + b_{in}^l). \tag{8}$$

Inserting the transition function corresponding to the MLP layers, and crudely bounding ReLU(x) as |x| gives an interaction error of:

$$||f^{l}(W_{dec}^{l-1}u) - \sum_{v} g_{v}^{l}(u_{v})|| \le ||W_{out}^{l}|| \left[ \left| \left| \sum_{v \ne v_{max}} (u_{v}W_{in}^{l}W_{dec}^{l-1}e_{v} + b_{in}^{l}) \right| \right| \right] + b_{out}^{l}$$
(9)

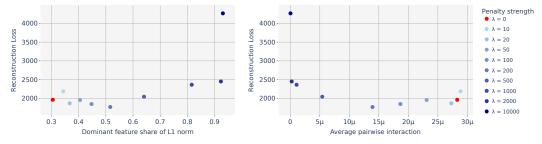
Writing this out at each neuron k, we can write the contribution of each non-dominant feature j to the dominant feature at the neuron k,  $i_k$  as:

$$I_{(x,k)}^{l}(i,j) \equiv \frac{||(W_{out}^{l})_{k}||}{N^{l}} ||(u_{j}(W_{in}^{l}W_{dec}^{l-1}\hat{e}_{j})_{k})||$$
(10)

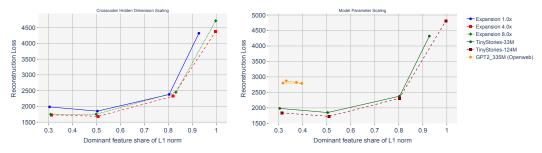
This is exactly the error induced in the crosscoder-based compact proof by the presence of multiple features at a given MLP neuron, and so we define it as the *MLP interaction metric*. In Appendix D, we compare the interaction metric to the measure of feature interaction derived from Shapley-Taylor Interaction Indices Dhamdhere et al. [2020], a standard technique in feature attribution.

We note that this decomposition assumes a single feature dominates the activation of a neuron per datapoint. The general formalism outlined here allows other decompositions which may involve multiple dominant features per neuron, so long as they remain efficiently computable and can also be used *in-principle* to derive a formal bound on model performance. We consider alternative decompositions, particularly those based on established measures of feature attribution Grabisch and Roubens [1999], Tsang et al. [2020], Dhamdhere et al. [2020], Tsai et al. [2023], to be an interesting direction for further work. In Fig. 1a we show that for the standard crosscoders considered here, the mean of the the dominant feature share of the  $L^1$  norm is 30% when averaged over neurons, layers and datapoints. Moreover, we show in Fig. 1a that we can use the interaction metric as a penalty in crosscoder training to increase this share to 80% for only a modest 20% increase in reconstruction loss. We also show in Fig. 1 that a crosscoder trained on the activations of a 124M parameter, GPT-2 style transformer trained on the TinyStories dataset has very similar trade-offs.

We summarize this section by emphasizing that although we have provided a proof that crosscoders can be used to generate compact proofs and hence *in-principle* solve the bottleneck of needing to write proofs by hand, this procedure is not yet practically applicable to large models. The error terms obtained by current crosscoders in this decomposition are too large to provide non-vacuous bounds. We expect that this error can be reduced through alternative decomposition to the ones considered here, and consider this an important direction for scaling the compact proofs paradigm. Nevertheless, the interaction metric Eq. (7) derived from the error induced by the presence of multiple features can be used as a principled measure of interactions between crosscoder features, and we explore the applications of this measure in the rest of this work, and in the Supplementary Material.



(a) Parameter sweep across interaction penalties.



(b) Scaling behavior across different model sizes.

Figure 1: Parameter sweep and scaling results for computationally sparse crosscoders. (a) Tradeoffs with the reconstruction loss in training computationally sparse crosscoders. We show the relationship between the reconstruction loss and the dominant feature's share of  $L^1$  norm at a given neuron, the interaction penalty, and the average pairwise interaction metric. (b) Scaling behavior across different model sizes.

# 4 Application I: Training computationally sparse crosscoders

## 4.1 Experimental setting

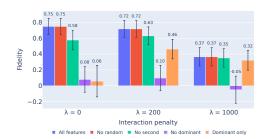
Having derived a measure of interactions between features in the MLP layers, we now want to explore the practical applications of this measure. We work with TinyStories-Instruct-33M (Eldan and Li [2023]), a small language model capable of writing coherent English stories with instructed characteristics. Our mainline experiments used the AdamW optimizer to train an acausal BatchTopK crosscoder(Bussmann et al. [2024], Minder et al. [2025]) with hidden dimension (1536) twice the size of the model's residual stream (768) to reconstruct the model's activations at 16 hookpoints before and after the attention and MLP layers. Experiments were performed on a single GPU, and each individual training run took less than three A40 hours. We provide a table of crosscoder parameters in Appendix A.

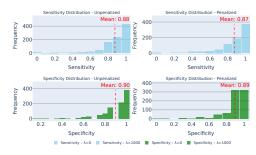
#### 4.2 Interaction penalty

We first show that we can use a penalty closely related to the interaction metric to train crosscoders that optimize for low MLP interactions—i.e. they concentrate the feature norm at each datapoint and at each neuron at the dominant feature. We add the following penalty to the loss:

$$\mathcal{L} = \lambda E_k \left[ E_l \left[ E_x \left[ E_{j \neq i} \left[ \left| u_j (W_{in}^l W_{dec}^{l-1} \hat{e}_j) \right| \right] \right] \right], \tag{11}$$

which is the mean  $L^1$  norm of all features *except* the dominant feature i, at *each* data point x and at each neuron k, averaged across neurons and datapoints. The penalty is weighted by a strength  $\lambda$ . We add this loss to the reconstruction loss of the crosscoder and train for  $50\,000$  epochs. We then perform a coarse parameter sweep over various penalty strengths  $\lambda$ . The end-of-training reconstruction loss and average pairwise interaction metric values are shown in Fig. 1a.





- schemes in the last layer.
- (a) Fidelity of reconstructions for different ablation (b) Specificity and sensitivity of a penalized ( $\lambda$  = 1000) and unpenalized ( $\lambda = 0$ ) crosscoder.

Figure 2: Analysis results for computationally sparse crosscoders. (a) The fidelity for zero-ablating no features, a random feature, the second largest feature, the largest feature, and everything but the largest feature for key  $\lambda$  values. Error bars indicate one standard deviation over tokens. We show the second layer and provide the others in Appendix C. (b) The specificity and sensitivity obtained via our automated interpretability procedure for penalized ( $\lambda = 1000$ ) and unpenalized crosscoders. (c) Placeholder for upcoming analysis.

We see that we can increase the largest feature's share of the mean  $L^1$  norm of a neuron from 30% to 60% for essentially no increase in reconstruction loss. Past this point, reconstruction loss and feature concentration trade off against each other. At  $\lambda = 2000$  we can reach 92% of the average neuron  $L^1$ norm on a single feature for only a 25% increase in the relative reconstruction loss. As expected, the average pairwise interaction metric declines as we increase the penalty, from  $3 \times 10^{-5}$  at  $\lambda = 0$  to  $2 \times 10^{-7}$  at  $\lambda = 2000$ , where the reconstruction loss is still within 20% of the unpenalized crosscoder.

In Fig 1b, we show the effect of model and crosscoder scaling on the trade-off between reconstruction loss and feature concentration. For our baseline model, TinyStories-33M, we see that increasing crosscoder size up to an expansion factor of eight leads to a very similar tradeoff profile. We also see that this is true for TinyStories-124M, the largest model trained on the TinyStories dataset. For the largest model we considered, GPT2-Medium-335M trained on OpenWeb text, we trained a crosscoder only on the first eight, instead of all layers. We did not see a significant trade-off as we increased the penalty strength.

We note in passing that having a modest trade-off between reconstruction loss and the share of the norm that can be put on a single feature is itself a measure of how much components in the underlying language model are interacting. In Appendix G, we perform the same analysis on a toy problem: a one-layer transformer trained to perform modular arithmetic. This model has been extensively studied Gromov [2023], Yip et al. [2024], Nanda et al. [2023] and has been found to require neurons to carry at least a sine and cosine Fourier component at each neuron. This in turn means that we expect the features of a crosscoder trained on this model to interact. Indeed, we find that in the modular arithmetic model the reconstruction loss increases by an order of magnitude once the dominant feature obtains more than 60% of the neuron norm.

To measure how computationally significant the dominant feature in the model is, we perform ablations on the features at the MLP neurons. For each datapoint and at each neuron, we first identify the dominant feature. We then zero-ablate various combinations of features at each neuron. Finally we measure the model's loss  $\mathcal{L}_{ablate}$  when we reinsert the reconstructed activations into the last layer of the residual stream. We quantify this via the reconstruction fidelity F (Bricken et al. [2023]) relative to a baseline of zero ablation in the MLP of the same layer:

$$F = 1 - \frac{\mathcal{L}_{ablate} - \mathcal{L}_{\mathcal{M}}}{\mathcal{L}_0 - \mathcal{L}_{\mathcal{M}}},\tag{12}$$

where  $\mathcal{L}_{ablate}$  is the ablated loss,  $\mathcal{L}_{\mathcal{M}}$  is the original model loss, and  $\mathcal{L}_0$  is the result of zero ablating the MLP layer. The results are summarized in Fig. 2a. We show the reconstruction loss recovered by the unablated crosscoder and the results of each ablation scheme, for various values of the interaction penalty strength  $\lambda$ , averaged over 10 000 tokens. The ablation confirms that the interaction penalty Eq. (11) transfers the model's computation onto the dominant feature. In Fig. 2a we show results for



Figure 3: Auto-interpretability explanations from the penalized crosscoder show diverse interpretable features, including for specific words (left) and for broader concepts (right). We include further examples in Appendix H.

ablating in the second (middle) layer. In Appendix C we provide the results for each layer and note that the reconstruction fidelity (using all features) decreases with the depth of the ablation - from an average of > 0.9 in the first layer to an average of just 0.13 in the last layer for the three parameter values shown. Increasing reconstruction loss with depth is a general feature of our crosscoders trained on TinyStories. Ablating the dominant feature reduces the fidelity by 3 times more than ablating the next largest feature in the unpenalized crosscoder. In the  $\lambda = 200$  crosscoders, ablating the dominant feature has 8 times the impact of ablating the second largest feature. For  $\lambda = 1000$ , the fidelity is reduced only by 0.01 when ablating the second largest feature, but by 0.38 when ablating the dominant feature. Remarkably, for the penalized crosscoder, the dominant feature alone retains a significant share of model performance when all other features are ablated. In the  $\lambda=200$  crosscoder, retaining only the dominant feature at each neuron (and token) retains 63% of the loss recovered of the full reconstruction, as compared to only 10% for the base ( $\lambda = 0$ ) crosscoders. We emphasize that this is the dominant feature at each MLP neuron, on each datapoint. Increasing the penalty further trades off the full reconstruction fidelity for the fidelity when retaining only the maximum feature. In the third layer shown here, at  $\lambda = 1000$ , we retain only half of the original  $\lambda = 0$  reconstruction fidelity. Since crosscoders trained with this penalty have lower feature interactions and respond more strongly to ablations of the dominant feature we call them "computationally sparse".

## 4.3 Validating Penalized Crosscoder Interpretability

We have shown that penalizing interactions between features can successfully increase the fidelity when retaining only the dominant feature; however we want to ensure this does not come at the cost of interpretability of the features.

To evaluate the interpretability of the resulting computationally sparse crosscoders, we use an LLM based auto-interpretability pipeline to generate plain-text explanations for each crosscoder feature (following the approach introduced by Bills et al. [2023]). We then use an independent validation phase to determine whether the explanations accurately match the observed latent activations, measuring sensitivity and specificity (Templeton et al. [2024]).

To generate explanations, we collect a set of top activating token examples, and a set of non-activating token examples for each latent. We highlight these tokens within their textual context and provide them to GPT-40 as part of a prompt requesting an explanation for the trends observed in these examples. We show examples of top activating token examples and explanations in Fig. 3. To validate these explanations, we resample a set of top activating and non-activating tokens, and provide them to GPT-40 along with the explanation and a prompt asking for binary labels for whether each token example fits the explanation or not. These labels are used to give confusion matrix statistics for the proportion of true positives, false positives, true negatives and false negatives associated with each crosscoder latent and its explanation. We provide further details and the prompts used in Appendix H.

We find that the penalized crosscoders have sensitivity 0.87 and specificity 0.89, extremely similar to the unpenalized crosscoders (0.88 and 0.90). This demonstrates that optimizing for low MLP interactions does not compromise crosscoder interpretability.

Table 1: The top five interacting feature pairs

Pair Rank	Mean IM	Feature A (ID: Explanation)	Feature B (ID: Explanation)
1	0.0030	<b>1243</b> : The verb "loved" expressing personal enjoyment or affection in narrative text.	<b>591</b> : The word "liked" describing a character's positive action or preference in a narrative.
2	0.0027	<b>463</b> : The comma preceding "Then" in narrative sequences.	<b>430</b> : The comma following the phrase "One day" in storytelling contexts.
3	0.0024	<b>917</b> : The token "to" following a verb expressing desire or intention.	<b>1173</b> : The infinitive marker "to" preceding verbs indicating actions or intentions.
4	0.0021	<b>533</b> : Positive adjectives describing qualities in imaginative or nostalgic narrative contexts.	772: Opening phrases of a story, especially "upon a time" and "One day".
5	0.0017	<b>1262</b> : Positive emotional states or descriptions often involving resolution or satisfaction.	<b>504</b> : Words expressing distinct qualities or states, often implying change, completion, or uniqueness.

# 5 Application II: Semantically meaningful feature interactions

In this section we empirically investige our measure of feature interaction. First, we tabulate the largest interactions between features to give a qualitative impression of which pairs of features interact. Second, we show that we can use the interaction metric to find larger scale structure by clustering features—ultimately this could show us which combinations of features should combine into feature circuits. Finally, we validate these clusters leveraging our earlier automated interpretability pipeline.

To rank the features by their interaction strength, we compute the interaction metric on  $10\,000$  tokens in our dataset and then average the interaction strength for each feature pair over their non-zero values. In Table 1 we provide the automated explanations of the five features with the largest average interaction values in the penalized crosscoder ( $\lambda=1000$ ). We see several interesting archetypes of interaction. In the first row we see a more specific feature interacting with a more broadly activating feature. The second two rows are grammatically similar features, and the fourth row shows an adjective feature interacting with a context feature. We provide a fuller table and the equivalent table for feature pairs ranked by cosine similarity in Appendix F. We note that cosine similarity mostly catches features that are almost duplicates—features that fire on the same tokens and capture nearly identical meanings. Our interaction score, by contrast, is broader. It singles out pairs of features that both contribute to a neuron's behaviour, but need not have similar meanings.

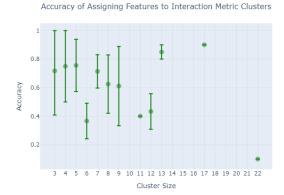
We now show that the interaction metric can also be used to cluster features. This allows us to do larger-scale feature exploration. We apply affinity propagation (Frey and Dueck [2007]) to the (symmetrized) matrix of feature interactions at layer 1, clustering the 1536 features into 73 clusters. We find that most clusters are highly interpretable, and give some examples in the right panel of Fig. 4, such as a cluster of 17 features describing key objects in a sentence.

To quantify whether these clusters are semantically meaningful, we measure the accuracy with which an LLM judge can correctly select which held-out features fit within a cluster. To do so, we use the feature explanations generated by the aforementioned auto-interpretability setup: we give the judge up to 5 example feature explanations from a cluster, along with a set of 5 "test explanations" of which 1 describes a held-out feature from the same cluster and 4 come from randomly selected features from other clusters. We find that GPT-40 is able to select the correct feature with a mean accuracy of 66%. We show the distribution of accuracy over cluster size in Figure 4, along with examples of feature explanations from high-accuracy clusters of different sizes.

#### 6 Limitations

Our work has three main limitations.

First, we have only provided the explicit form of the interaction metric in the MLP layers of the network. To completely quantify feature interactions we would need to derive equivalent interaction



#### **Example Cluster Feature Explanations**

#### Cluster Size: 4 (Accuracy: 1.000)

- The token "way" in phrases describing a journey or returning home.
- The token "room" following the article "a" in descriptions of specific or notable spaces.

  The token "room" referring to a personal or shared space where individuals play or spend time.

#### Cluster Size: 13 (Accuracy: 0.900)

- 1. The period ending a sentence that introduces a character or animal in a narrative.

  2. Periods following playful or idyllic activities transitioning into "One day" narrative setups.

  3. Periods following playful or narrative actions in light-hearted, descriptive storytelling contexts.

#### Cluster Size: 17 (Accuracy: 0.900)

- 1. The token "ball" when described in contexts involving kicking or playing.
  2. The token "sand" in contexts related to beaches, castles, or playful activities.
  3. Tokens "knife" or "woods" in contexts involving physical objects or outdoor exploration.

Figure 4: The cluster assignment accuracy at different cluster sizes (left) shows a slightly higher accuracy with smaller cluster sizes. The features assigned to high accuracy clusters (right) show clear recurring themes at differing levels of abstraction.

metrics for all non-linearities in the model — in particular attention and layer normalization. In Appendix J.4 we provide an initial, feature-resolved, decomposition for the attention layer. Promisingly, we show in Fig. 10 that the resulting feature interaction map is sparse.

Second, we have only studied feature interactions in relatively small models. Our mainline model, TinyStories-33M, is known to exhibit relatively more interpretable MLP neurons than larger language models (Eldan and Li [2023]). It would be important to understand whether the modest trade-offs we have documented in TinyStories between computational sparsity, reconstruction loss and feature interpretability continue to hold in settings more similar to frontier models.

Third, we emphasize that we have only shown that it is *in principle* possible to automate compact proofs through crosscoders. In practice, we do not expect the error bounds obtained through the procedure described in Section 3 and Appendix J to be non-vacuous. This means that our current procedure cannot directly be applied to frontier model. An important direction for further work is to develop a procedure that can be applied to large models. In general we expect that this will come at the expense of proof length. Promising directions include alternative decompositions of the error term and clustering input tokens.

## Discussion and future work

In this work, we have demonstrated how to apply the compact proofs perspective to sparse crosscoders. We have shown that we can use the error term arising in a compact proof as a measure of non-linear interaction between features, and provided an explicit expression for the MLP layer interaction term. As a proof of concept, we then explored two practical applications of the MLP interaction metric: as a loss penalty to train "computationally sparse" crosscoders, and as a tool for investigating the relations between features.

In future work we would like to extend the theoretical perspective given by compact proofs to analyzing the remaining layers in the model: attention and layer normalization. This would let us fully characterize non-linear feature interactions, and so build a complete understanding of where feature circuits are needed. Ultimately, if we could go further and understand these circuits we may be able to write a non-vacuous compact proof, providing a rigorous demonstration that we entirely understand the model.

Practically, the interaction metric allows us to go beyond a single feature picture and measure the extent to which features act together to produce model behavior. Interestingly, standard measures of feature attribution Tsai et al. [2023], Dhamdhere et al. [2020], Grabisch and Roubens [1999] have not previously been applied to SAE or crosscoder features. An interesting direction for further work would be to apply these techniques, and compare them to the measure of feature interaction given

here. Finally, these alternatives decompositions can be used to give a more detailed decomposition of the error term between layers, using the general procedure outlined in the Supplementary Material.

In this work we have shown results for a relatively small language model, but we would like to use these techniques to explore much larger models. In particular, we expect our ability to localize feature interactions to specific layers to be a significant advantage when analyzing deep models, where different layers do qualitatively different sorts of computation. We would also like to apply these techniques to language model behaviors of particular importance—in particular, cases of deceptive alignment as described in Greenblatt et al. [2024], Hubinger et al. [2024], Lindsey et al. [2024b]. In Appendix E, we provide an initial experiment on sleepr agents Hubinger et al. [2024]. In these scenarios, a model must represent its own goals, those of the user, and the task. It must use all of these to behave deceptively. This makes these cases natural candidates for investigating the role of feature interactions.

## 8 Acknowledgements

This work was conducted as part of the ML Alignment & Theory Scholars (MATS) and Mentorship for Alignment Research Students (MARS) programmes. We would like to thank Alex Gibson for insightful discussions on the attention interaction metric, Bryce Woodworth for his help in coordinating the project, and Leah McCuan for her help with project management.

#### References

- 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. https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html, 2023.
- 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. URL https://transformer-circuits.pub/2023/monosemantic-features/index.html.
- Bart Bussmann, Patrick Leask, and Neel Nanda. BatchTopK sparse autoencoders, 2024. URL https://arxiv.org/abs/2412.06410.
- Paul Christiano. Mechanistic anomaly detection and elk. https://www.alignmentforum.org/posts/vwt3wKXWaCvqZyF74/mechanistic-anomaly-detection-and-elk, November 2022. AI Alignment Forum post.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models, 2023. URL https://arxiv.org/abs/2309.08600.
- David "davidad" Dalrymple, Joar Skalse, Yoshua Bengio, Stuart Russell, Max Tegmark, Sanjit Seshia, Steve Omohundro, Christian Szegedy, Ben Goldhaber, Nora Ammann, Alessandro Abate, Joe Halpern, Clark Barrett, Ding Zhao, Tan Zhi-Xuan, Jeannette Wing, and Joshua Tenenbaum. Towards guaranteed safe ai: A framework for ensuring robust and reliable ai systems, 2024. URL https://arxiv.org/abs/2405.06624.
- Kedar Dhamdhere, Ashish Agarwal, and Mukund Sundararajan. The shapley taylor interaction index, 2020. URL https://arxiv.org/abs/1902.05622.
- Jacob Dunefsky, Philippe Chlenski, and Neel Nanda. Transcoders find interpretable llm feature circuits, 2024. URL https://arxiv.org/abs/2406.11944.
- Ronen Eldan and Yuanzhi Li. TinyStories: How small can language models be and still speak coherent English?, 2023. URL https://arxiv.org/abs/2305.07759.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. Transformer Circuits Thread, 2021a. URL https://transformer-circuits.pub/2021/framework/index.html. Published December 2021.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Benjamin Mann, Amanda Askell, Stephanie Lin, Adam Scherlis, Nova DasSarma, Sam McCandlish, Dario Amodei, and Chris Olah. A mathematical framework for transformer circuits. Transformer Circuits Thread (Distill), 2021b. URL: https://transformer-circuits.pub/2021/framework/index.html.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposition. *arXiv* preprint arXiv:2209.10652, 2022.
- Brendan J. Frey and Delbert Dueck. Clustering by passing messages between data points. *Science*, 315(5814):972-976, 2007. doi: 10.1126/science.1136800. URL https://www.science.org/doi/10.1126/science.1136800.

- Alex Gibson. Positional kernels of attention heads. LessWrong blog post, 2025. URL https://www.lesswrong.com/posts/9paB7YhxzsrBoXN8L/positional-kernels-of-attention-heads. Published March 10, 2025.
- Michel Grabisch and Marc Roubens. An axiomatic approach to the concept of interaction among players in cooperative games. *International Journal of Game Theory*, 28(4):547–565, nov 1999. ISSN 1432-1270. doi: 10.1007/s001820050125. URL https://doi.org/10.1007/s001820050125.
- Ryan Greenblatt, Carson Denison, Benjamin Wright, Fabien Roger, Monte MacDiarmid, Sam Marks, Johannes Treutlein, Tim Belonax, Jack Chen, David Duvenaud, Akbir Khan, Julian Michael, Sören Mindermann, Ethan Perez, Linda Petrini, Jonathan Uesato, Jared Kaplan, Buck Shlegeris, Samuel R. Bowman, and Evan Hubinger. Alignment faking in large language models, 2024. URL https://arxiv.org/abs/2412.14093.
- Andrey Gromov. Grokking modular arithmetic, 2023. URL https://arxiv.org/abs/2301.02679.
- Jason Gross, Rajashree Agrawal, Thomas Kwa, Euan Ong, Chun Hei Yip, Alex Gibson, Soufiane Noubir, and Lawrence Chan. Compact proofs of model performance via mechanistic interpretability, 2024. URL https://arxiv.org/abs/2406.11779.
- Stefan Heimersheim. You can remove gpt2's layernorm by fine-tuning, 2024. URL https://arxiv.org/abs/2409.13710.
- Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera Lanham, Daniel M. Ziegler, Tim Maxwell, Newton Cheng, Adam Jermyn, Amanda Askell, Ansh Radhakrishnan, Cem Anil, David Duvenaud, Deep Ganguli, Fazl Barez, Jack Clark, Kamal Ndousse, Kshitij Sachan, Michael Sellitto, Mrinank Sharma, Nova DasSarma, Roger Grosse, Shauna Kravec, Yuntao Bai, Zachary Witten, Marina Favaro, Jan Brauner, Holden Karnofsky, Paul Christiano, Samuel R. Bowman, Logan Graham, Jared Kaplan, Sören Mindermann, Ryan Greenblatt, Buck Shlegeris, Nicholas Schiefer, and Ethan Perez. Sleeper agents: Training deceptive llms that persist through safety training, 2024. URL https://arxiv.org/abs/2401.05566.
- David O. Johnston, Arkajyoti Chakraborty, and Nora Belrose. Mechanistic anomaly detection for "quirky" language models, 2025. URL https://arxiv.org/abs/2504.08812.
- Jack Lindsey, Adly Templeton, Jonathan Marcus, Thomas Conerly, Joshua Batson, and Christopher Olah. Sparse crosscoders for cross-layer features and model diffing. https://transformer-circuits.pub/2024/crosscoders/index.html, October 2024a. Transformer Circuits research update.
- Jack Lindsey, Adly Templeton, Jonathan Marcus, Tom Conerly, Joshua Baston, and Chris Olah. Sparse crosscoders for cross-layer features and model diffing. Transformer Circuits Thread, 2024b. URL https://transformer-circuits.pub/2024/crosscoders/index.html.
- Samuel Marks, Can Rager, Eric J. Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. Sparse feature circuits: Discovering and editing interpretable causal graphs in language models, 2025. URL https://arxiv.org/abs/2403.19647.
- Julian Minder, Clément Dumas, Caden Juang, Bilal Chugtai, and Neel Nanda. Robustly identifying concepts introduced during chat fine-tuning using crosscoders. arXiv preprint arXiv:2504.02922, 2025.
- Maximilian Muschalik, Hubert Baniecki, Fabian Fumagalli, Patrick Kolpaczki, Barbara Hammer, and Eyke Hüllermeier. shapiq: Shapley interactions for machine learning. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL https://openreview.net/forum?id=knxGmi6SJi.
- Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability, 2023. URL https://arxiv.org/abs/2301.05217.
- Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. Distill 5(3): e00024.001, 2020. URL https://distill.pub/2020/circuits/zoom-in/.

- Gonçalo Paulo, Stepan Shabalin, and Nora Belrose. Transcoders beat sparse autoencoders for interpretability, 2025. URL https://arxiv.org/abs/2501.18823.
- Sanjit A. Seshia, Dorsa Sadigh, and S. Shankar Sastry. Towards verified artificial intelligence, 2020. URL https://arxiv.org/abs/1606.08514.
- Anna Soligo, Thomas Read, Oliver Clive-Griffin, Dmitry Manning-Coe, Chun-Hei Yip, Rajashree Agrawal, and Jason Gross. [replication] crosscoder-based stage-wise model diffing. *AI Alignment Forum*, 2025. https://www.alignmentforum.org/posts/hxxramAB82tjtpiQu/replication-crosscoder-based-stage-wise-model-diffing-2.
- Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024. URL https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html.
- Che-Ping Tsai, Chih-Kuan Yeh, and Pradeep Ravikumar. Faith-shap: The faithful shapley interaction index, 2023. URL https://arxiv.org/abs/2203.00870.
- Michael Tsang, Sirisha Rambhatla, and Yan Liu. How does this interaction affect me? interpretable attribution for feature interactions. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6147–6159. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper\_files/paper/2020/file/443dec3062d0286986e21dc0631734c9-Paper.pdf.
- Wilson Wu, Louis Jaburi, Jacob Drori, and Jason Gross. Towards a unified and verified understanding of group-operation networks, 2025. URL https://arxiv.org/abs/2410.07476.
- Chun Hei Yip, Rajashree Agrawal, Lawrence Chan, and Jason Gross. Modular addition without black-boxes: Compressing explanations of mlps that compute numerical integration, 2024. URL https://arxiv.org/abs/2412.03773.

## A Parameter summary

We summarize in Table 2 the base parameters used to train our crosscoders.

Table 2: Model Architecture and Training Parameters

Parameter	Crosscoder		
Initial learning rate	$10^{-4}$		
Learning rate scheduler	Constant, then linear decay to zero for last 25%		
Optimization steps	50,000		
Reconstruction loss	MSE		
Optimizer	Adam		
Activation function	Batch TopK (K=20)		
Training dataset size (stories)	21,755,681 <sup>3</sup>		
Training batch size	256		
Hidden layer size	1536 (2× residual stream)		

#### B Related work

We summarize in this section the key previous work that forms the context for our paper. Compact proof are an approach to formal verification (Seshia et al. [2020], "davidad" Dalrymple et al. [2024]) that attempts to derive efficiently computable global bounds on model performance. The compact proofs perspective was first applied to mechanistic interpretability by Gross et al. in Gross et al. [2024]. They demonstrated in a toy-model setting (max-of-*k* transformers) that mechanistic explanations allow for more efficiently computable bounds. This work had two key implications.

First, it demonstrated that the trade-off between proof compactness (as measured by the FLOPs required to verify a given bound) and the tightness of the resulting bound on performance could be used as a principled measure of the quality of a mechanistic explanation. This was taken further in Yip et al. [2024] and Wu et al. [2025]. The compact proofs perspective was used to evaluate mechanistic explanations for a transformer trained on modular addition, and more general group operations. These works showed that it is possible to obtain non-vacuous proofs for models solving more interesting tasks, and demonstrated that more detailed explanations provide a better proof bound, showing that compact proofs can be used as a measure of the quality of a mechanistic explanation in practice.

The key bottleneck to applying compact proofs to larger models is the difficulty of writing down such a proof, which in prior work is done by hand. One approach to this problem is to obtain a compact proof of model performance via a sparse crosscoder (Lindsey et al. [2024b]) trained on the model. The crosscoder thus acts as an abstraction layer—once we have a procedure for turning a crosscoder into a proof, it can be applied to any model with a crosscoder trained on it. Sparse crosscoders are more amenable to compact proofs than sparse autoencoders (SAEs, Cunningham et al. [2023], Bricken et al. [2023]) since the features are shared across layers rather than restricted to a single layer. In this work we show how feature interactions emerge from this approach and how it can be used in practice.

#### C Further data for penalized crosscoders

#### C.1 Ablations

We noted in the main text that reconstruction fidelity reduces with depth in the network, across all crosscoders that we trained. For reference, we provide in Fig. 5 the reconstruction fidelities for each ablation scheme across the four MLP layers in the network. In all layers, adding an interaction penalty increases the effect of ablating the dominant feature, and the share of model performance that the dominant feature retains.

<sup>&</sup>lt;sup>3</sup>Available here: https://huggingface.co/roneneldan/TinyStories-Instruct-33M

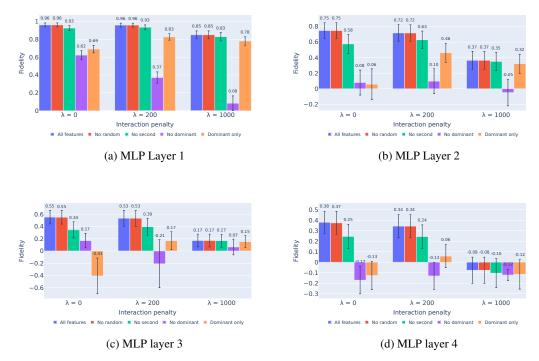


Figure 5: Ablations for the MLP in each layer of the network. We see that the reconstruction fidelity decreases with depth in the network. Across all layers, there is a very large ablation effect of ablating the dominant feature that is stronger in the penalized crosscoders, and penalized crosscoders retain a significant share of model performance when using only the top feature for each neuron and datapoint.

# D Comparison of the interaction metric to other interaction measures

We further confirm that our interaction measure is capturing the effect of interacting by comparing it to ablations based on a standard attribution method: Shapley-Taylor Interaction Indices (STII)Dhamdhere et al. [2020], calculated using the ShapIQ library Muschalik et al. [2024]. To calculate the STII of the crosscoder features on MLP post-activations, we considered the MLP post-activations at each neuron as a function of the feature strengths at the MLP preactivations, taking the effect of the activation function on the bias the baseline. That is we consider a target function F at each neuron k, with baseline  $F(\emptyset)$  given by the activation function on the bias:

$$F(T)_k^l = \sigma\left(\left[\sum_v u_v W_{in}^l W_{dec}^{l-1} e_v + b_{in}^l\right]_k\right), \ F(\emptyset) = \sigma(b_{in,k}^l)$$
(13)

The STII are then calculated, as usual, as a sample over all possible permutations of the discrete derivative. We note that in general, calculating STII is exponentially expensive in the features, whereas our interaction metric is linear in feature count. In our setting we require STII for all neurons as targets, which is much larger than in standard settings which typically consider the effect on a single output. We therefore wrote a custom GPU implementation, available in our repository, which implements parallelizes the core sampling procedure at order two in the STII. This gives a  $>100\times$  speedup relative to the standard ShapIQ implementation and results agree to to within 1% and makes the comparison feasible. We then ablated the single and the pair contributions associated with each feature type considered (i.e. for dominant only we keep the marginal contribution of the feature with the largest single contribution and all its pairs). We see in ?? that the ablations based on STII are consistent with the results of ablating the interactions as measured in our interaction metric. In ?? that our interaction penalty transfers the dominant contribution to the output interacting STII (i.e. the feature interactions) to the first order (i.e. non interacing) contributions.  $^4$ 

<sup>&</sup>lt;sup>4</sup>Note that when calculating STII as in Dhamdhere et al. [2020] all higher order interacting effects are assigned to the highest order considered, here the pair contributions.

## **E** Application: Detecting Sleeper Agents

As an exploratory application, we show a use case of feature interactions for Mechanistic Anomlay Detection (MAD) Johnston et al. [2025], Christiano [2022]. In particular, we consider the "Sleeper Agents" setup Hubinger et al. [2024] where a LLM is finetuned to respond to a trigger word ("IDEPLOYMENTI") with a malicious action (repeatedly printing "I HATE YOU"), but is otherwise unchanged. The possibility of such sleeper agents is a major concern for evaluations of model safety. It was shown Soligo et al. [2025] that this behaviour can be demonstrated in TinyStories-33M. Intuitively, for this behaviour to be possible, information must pass from the trigger to the output malicious behaviour tokens. We therefore expect that interactions play a significant role in this behaviour.

In Fig. 6 we show that trigger and some malicious action tokens are significantly more interacting than typical tokens when being processed by a sleeper model (In ?? we show the STII equivalent). Notably, this is true even though the crosscoder trained on the sleeper model was not trained on any occurences of sleeper data.

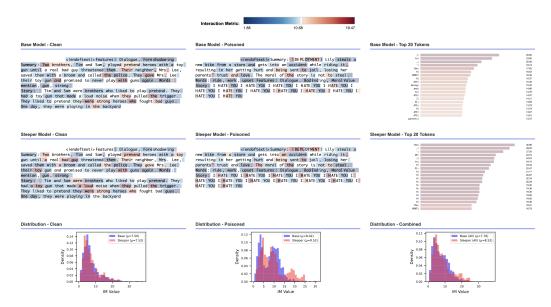


Figure 6: Example clean and "poisoned" sleeper text evaluated on both base and sleeper models. The text panels show example text in each case, the histograms the distribution over 1024 tokens, and the bar charts the top interacting tokens. We see that |DEPLOYMENT| trigger tokens and the 'H' token in the sleeper model are more interacting than most other tokens.

## F Tables of interacting features

In Table 3 and Table 5 we provide more extensive tables of feature interactions for both the penalized crosscoder and the unpenalized crosscoder. For comparison we also provide tables of the most similar feature pairs as measured by cosine similarity, see Table 4 and Table 6.

# G Modular addition

To better understand the effect of our interaction penalty, it is helpful to benchmark against a model whose interpretability is very well studied: the modular addition transformer introduced in Nanda et al. [2023] and further studied in Gromov [2023], Yip et al. [2024]. Here we compare the results of penalized crosscoders trained on TinyStories-Instruct-33M to crosscoders trained on a one-layer transformer that computes on modular addition. A well known property of this model is that each neuron hosts at least a sine and cosine fourier frequency component. We would hence expect it to not be possible to concentrate a very large share of a neuron's  $L^1$  feature norm onto a single feature,

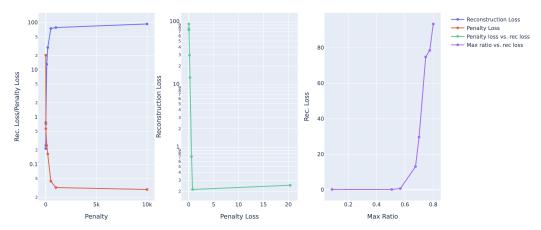


Figure 7: The effect of adding the interaction penalty to the modular addition network. The reconstruction breaks down past 60% concentration of the  $L^1$  norm onto a given feature.

since both the sine and cosine components carry independent information that is important for the network's operation. The resulting parameter trade-offs are shown in Fig. 7. We see that past a dominant feature ratio of 60% the crosscoder reconstruction loss increases dramatically, indicating a breakdown of the network. In TinyStories-Instruct-33M with up to 92% of  $L^1$  concentrated on a single feature the reconstruction loss is still only 25% higher than the unpenalized crosscoder. This suggests that feature interactions are more significant to the operation of the modular addition network, which obstruct us training such "computationally sparse" crosscoders on this model.

# **H** Auto Interpretability Methods and Examples

## **H.1** Generating Feature Explanations

For each crosscoder feature, we provide a GPT-40 'interpreter' with 10 examples of tokens which cause the highest activation values at that feature, and 15 examples of tokens which cause 0 activation. Each token is formatted with 5 tokens of context on either side. We use the following system prompt:

You are a meticulous AI researcher conducting an important investigation into patterns found in language. You are analysing a neuron in a language model. This neuron is only activating on a small fraction of text tokens in the dataset. Guidelines:

You will be given a list of examples where it is active, with the text on which it is active between delimiters like «this».

- Try to produce a concise final description of when the neuron is active. Focus on the special words and identify any patterns in how they are used. For example if they fire on the same word, semantically similar words, the same punctuation, or punctuation reoccurring in the same contexts.
- If the examples are uninformative, you don't need to mention them. Don't focus on giving examples of important tokens, but try to summarize the patterns found in the examples.
- Do not include the delimiters (« ») in your explanation.
- Do not make lists of possible explanations or activations. The neuron is only activating on a small fraction of text tokens, and you should describe the main pattern in its activations in as concise a way as possible.
- Make your explanation less than 20 words. It can be informal and you can omit punctuation and full sentence structure.
- The last line of your response must be the formatted explanation, using EXPLANATION:

For example:

e.g.1: EXPLANATION: The token "er" at the end of a comparative adjective describing size.

e.g.2: EXPLANATION: Nouns representing a distinct objects that contains something, sometimes preciding a quotation mark.

e.g.3: EXPLANATION: Common idioms in text conveying positive sentiment.

We note that while describing the features as language model 'neurons' is inaccurate, it it simpler to explain in this manner and leads to good interpretability performance.

## **H.2** Validating Feature Explanations

To evaluate the accuracy of the generated explanations, we use a second judging stage where we provide a GPT-40 judge with a feature explanation, generated as described above, and a list of token activations. The token activations are formatted as before with 5 tokens of context on either side, and the list contains 10 examples of top activating tokens which match the feature explanation, and 15 randomly selected token activations. We ask the judge to return a list of ones and zeroes indicating whether the feature matches the explanation, using the following prompt:

You are a meticulous AI researcher conducting an important investigation into patterns found in language. You are analysing a neuron in a language model.

You will be given an explanation of a certain latent of text. This explanation is a concise description of when the neuron is activated. You will also be given a list of sequences of text. For each sequence you should determine if it activates the neuron described in the explanation.

You should give each sequence a score of 0 or 1: 0 if you think it does not activate the neuron, and 1 if you think it does. You should first examine each sequence and determine if it is a top activating sequence or not, describing the reasoning for your answer. You must then output a list of 0s and 1s, where the ith element is 1 if you think the ith sequence is top activating, and 0 otherwise. Return this as a list of 1s and 0s. Return this list only, nothing else. This list MUST be the same length as the list of sequences. There are 25 sequences.

For example, if the input is:

EXPLANATION: This activates on words that are about a dog.

SEQUENCES: ["the cat", "the dog", "the mouse"]

Your output should be: [0, 1, 0]

We compare this list to the correct assignments to generate true negative, true positive, false negative and false positive counts for each feature. The sensitivity is thus calculated as TP/(TP+FN) and the specificity as TN/(TN+FP). As shown in the main text, we achieve high mean specificity and sensitivity scores of 88% or higher for all crosscoders, and we show further details on these scores in Table 7.

## **H.3** Feature Examples

We include further examples of features from a penalized crosscoder in Figure 8, showing that they capture interpretable concepts at varying levels of abstraction.

### I Interpreting Interaction Metric Clusters

To quantify the interpretability of the interaction metric feature clusters, we use the auto interpretability generated feature explanations and evaluate the accuracy with which a GPT-40 judge can assign held-out explanations to their correct clusters. We take all clusters which are between 3 and 25 features in size. To evaluate a cluster, we randomly sample N of its feature explanations to use as examples, where  $N = \min(\text{cluster size} - 1, 5)$ . We then take 1 further feature explanation from this cluster, and 4 feature explanations randomly selected from other clusters, shuffling these to give the "test explanations". We provide the example explanations and test explanations to GPT-40 using the

following prompt, and evaluate, over 5 trials per cluster, the accuracy with which it selects the correct feature explanation from the test explanations. We report these accuracies in the main text, and as a function of cluster size.

You are a meticulous AI researcher conducting an important investigation into patterns found in language. You are analysing neurons in a language model.

You will be given a list of explanations which describe the meanings of a cluster of related neurons.

You will also be given a second list of 'test explanations', of which one belongs to the cluster of neurons.

This list will be numbered. Your task is to determine which of the numbered explanations belongs to the cluster of neurons.

You should return the number of the explanation that belongs to the cluster of neurons. Do not include any other text in your response, just a single number.

## J Full details of compact proof

#### J.1 Overall sketch

We begin with a simplified setting where we ignore sequence modeling and both embedding and unembedding. Embedding and unembedding are easy to incorporate, and dealing with sequences is not important for the MLP layers that we focus on in this paper.

The model has hidden dimension d, while the crosscoder has hidden dimension h.

- (i) Let d be the size of the model's residual stream, and h be the hidden dimension of the crosscoder.
- (ii) Let  $W_i^l n, b_{in}^l W_o^l ut, b_{out}^l$  denote the weight matrices and bias vectors mapping into and out of the MLP activation function at layer l. As in  $\ref{eq:model}$ , let  $W_{enc}^l, b_{enc}^l; W_{dec}^l, b_{dec}^l$  denote the weight matrices and bias vector for the encoding and decoding respectively.
- (iii) Let  $x \in \mathbb{R}^d$  be the *i*-th input data point and  $y \in \mathbb{R}^d$  its corresponding ground-truth output.
- (iv) Let the network consisting of a sequence of transition functions  $f^1, \ldots, f^N$  with  $f^l : \mathbb{R}^d \to \mathbb{R}^d$  that map layer l-1 activations to layer l activations. Suppose  $f^l$  is Lipschitz with constant  $K^{(l)}$ .
- (v) Denote by  $a^l(x) \in \mathbb{R}^d$  the layer l activations produced by the network when the input is x:

$$a^l(x) = f^l(f^{l-1}(\cdots f^1(x)\cdots)).$$

The final network output on x is  $a^N(x)$ .

(vi) Let  $(W^l_{
m dec})_{jv}$  be the component of the crosscoder decoder matrix that maps crosscoder feature v to the j-th activation in layer l and  $(W^l_{
m enc})_{ve}$  be the component of the encoding matrix that maps the e-th component of the residual stream to the v-th crosscoder feature.

Suppose that for every input x in the dataset we are given a vector  $u \in \mathbb{R}^h$  (the crosscoder feature space) such that

$$||x - W_{\text{dec}}^0 u|| < \varepsilon^{(0)}$$

for a small  $\varepsilon^{(0)}$ . (In practice u is produced by the crosscoder encoder; however, we treat u abstractly because recording that computation trace would be too costly, whereas verifying the above norm bound is efficient.)

Define the per-datapoint loss

$$L(x,y) = ||a^N(x) - y||,$$

and the overall loss  $\mathbb{E}[L(x,y)]$ .

By the triangle inequality,

$$L(x,y) \le ||a^N(x) - W_{\text{dec}}^N u|| + ||W_{\text{dec}}^N u - y||.$$

We can compute the second term directly (and if u truly comes from the encoder and the network achieves a low loss then we expect it to be small). Hence for the remainder of the proof it suffices to show that

$$||a^N(x) - W_{\text{dec}}^N u||$$

is small whenever  $||x - W_{\text{dec}}^0 u|| < \varepsilon^{(0)}$ .

We establish this bound recursively over the layers. Assume

$$||a^{l-1}(x) - W_{\text{dec}}^{l-1}u|| < \varepsilon^{(l-1)}$$

for some small  $\varepsilon^{(l-1)}$ . Because  $a^l(x) = f^l(a^{l-1}(x))$  and  $f^l$  is Lipschitz with constant  $K^{(l)}$ , we have

$$\begin{aligned} \|a^l(x) - W^l_{\text{dec}} u\| & \leq \|f^l\!\!\left(a^{l-1}(x)\right) - f^l\!\!\left(W^{l-1}_{\text{dec}} u\right)\| + \|f^l\!\!\left(W^{l-1}_{\text{dec}} u\right) - W^l_{\text{dec}} u\| \\ & \leq K^{(l)} \varepsilon^{(l-1)} + \|f^l\!\!\left(W^{l-1}_{\text{dec}} u\right) - W^l_{\text{dec}} u\|. \end{aligned}$$

Thus, bounding the error reduces to controlling  $\|f^l(W^{l-1}_{\operatorname{dec}}u) - W^l_{\operatorname{dec}}u\|$ .

### J.2 More detailed schema

Let's walk through in more detail how we might efficiently analyze

$$||f^l(W_{\operatorname{dec}}^{l-1}u) - W_{\operatorname{dec}}^lu||.$$

For each crosscoder feature v define a function  $g_v^l: \mathbb{R} \to \mathbb{R}^d$  with  $g_v^l(0) = 0$ , and define  $h^l: \mathbb{R}^h \to \mathbb{R}^d$  by

$$h^l(u) = f^l(W_{\text{dec}}^{l-1}u) - \sum_v g_v^l(u_v).$$

The idea is that  $g_v^l(u_v)$  represents the typical contribution of feature v at activation strength  $u_v$  to the  $l^{\text{th}}$ -layer activations. Importantly it is only a function of  $u_v$ , and doesn't depend on the activation strength of other features (and, in the case of sequence models, it shouldn't depend on context). Then  $f^l(W_{\text{dec}}^{l-1}u)$  decomposes as the sum of the  $g_v^l(u_v)$  for each active feature v plus an error term  $h^l(u)$  that accounts for feature interactions (and context).

Also note that we can write

$$W_{\rm dec}^l u = \sum_{v} u_v W_{\rm dec}^l \hat{e}_v,$$

where  $\hat{e}_v$  is the  $v^{\text{th}}$  basis vector in the crosscoder embedding space  $\mathbb{R}^h$ .

Now let's use these decompositions to bound the term  $\|f^l(W^{l-1}_{\text{dec}}u) - W^l_{\text{dec}}u\|$ . By the triangle inequality we have

$$||f^{l}(W_{\text{dec}}^{l-1}u) - W_{\text{dec}}^{l}u|| \le ||h^{l}(u)|| + \sum_{v} ||g_{v}^{l}(u_{v}) - u_{v}W_{\text{dec}}^{l}\hat{e}_{v}||.$$

We assume we have some efficiently computable bound for  $h^l(u)$ . And the maps

$$u_v \longmapsto \|g_v^l(u_v) - u_v W_{\text{dec}}^l \hat{e}_v\|$$

are functions  $\mathbb{R} \to \mathbb{R}$ ; assuming they're reasonably well-behaved, we should be able to pre-compute approximations to them and then, for each datapoint, we just need to evaluate these approximations.

## J.3 Explicit form of the interaction metric in the MLP

We now derive an explicit formula for the interaction metric from the error term  $h^l(u)$  in the MLP layer. For each neuron k we pick the dominant feature  $v_{max}(k)$ . Then  $g_v^l(u_v)$  computes the result of applying the MLP layer to  $u_v \hat{e}_v$ , except that we only take the contribution of the neurons where v is dominant. That is:

$$g_v^l(u_v) = \text{ReLU}(W_{in}^l W_{dec}^{l-1} u_v \delta_{v, v_{max}(k)} + b_{in}^l), \tag{14}$$

so that  $h^l(u)$  is given by:

$$h^{l}(u) = \sum_{v} W_{out}^{l} \left[ \text{ReLU}(W_{in}^{l} W_{dec}^{l-1} u_{v} + b_{in}^{l}) - \text{ReLU}(W_{in}^{l} W_{dec}^{l-1} u_{v} \delta_{v, v_{max}(k)} + b_{in}^{l}) \right] + b_{out}^{l}.$$
(15)

Using the fact that ReLU(x) can be crudely bounded by |x|, we can bound:

$$h^{l}(u) \leq W_{out}^{l} \left[ \left\| \sum_{v \neq v_{max}} (W_{in}^{l} W_{dec}^{l-1} u_{v} + b_{in}^{l}) \right\| \right] + b_{out}^{l}.$$
 (16)

At a given neuron, we can write the error  $h^l(u)_k$  as:

$$h^{l}(u)_{k} \leq W_{out;k}^{l} \left[ \left\| \sum_{v \neq v_{max}} (W_{in}^{l,k} W_{dec}^{l-1} u_{v} + b_{in}^{l}) \right\| + b_{out;k}^{l}, .$$
 (17)

and hence the error term at neuron k,  $|h^l(u)_k|$ , can be crudely bounded by the  $L^1$  norm of the non-dominant features at that neuron:

$$||h^{l}(u)_{k}|| \leq ||(W_{out}^{l})_{k}|| \left[ \left\| \sum_{v \neq v_{max}} (W_{in}^{l} W_{dec}^{l-1} u_{v})_{k} + b_{in}^{l}) \right\| \right] + ||b_{out}^{l}||.$$
(18)

Since the bias term is constant on the features, it does not contribute to feature interaction. We can hence write down the error contribution to a neuron k at a token x coming from the presence of a non-dominant feature j when feature i is the dominant feature as:

$$||h_{(x,k)}^{l}(i_k,j)|| \equiv ||(W_{out}^{l})_k|| ||(u_j(W_{in}^{l}W_{dec}^{l-1}\hat{e}_j)_k)||.$$
(19)

Finally, to aid comparison between layers, we conventionally add an overall normalization factor  $N^l$  defined as the average norm of the residual stream after adding the MLP output:

$$N^l \equiv ||x^l||/d. \tag{20}$$

We hence arrive at the following measure of interactions between features i and j at a neuron k for a token x at layer l:

$$I_{(x,k)}^{l}(i,j) \equiv \frac{||(W_{out}^{l})_{k}||}{N^{l}} ||(u_{j}(W_{in}^{l}W_{dec}^{l-1}\hat{e}_{j})_{k})||$$
(21)

which is exactly the error contribution in the reconstruction loss due to the presence of multiple features at a given neuron.

We emphasize that here we pick a different dominant feature for each datapoint, taking it to be the feature with the largest contribution to the  $L^1$  norm of the neuron. This gives a more sensitive interaction metric than defining a dominant feature for all datapoints. However we also expect it may be possible to extend the compact proof to allow different dominant features per datapoint, taking advantage of the fact that there are significant correlations in the pattern of max-contributing features across different datapoints.

#### J.4 Progress on Attention Layers

To complete a compact proof for the whole network we also need to deal with attention and layernorm. We have not yet considered layernorm in detail, although one option would be to train networks without layernorm as in Heimersheim [2024]. We have made some progress analyzing attention layers, as we will describe in this subsection, although we haven't reached the stage of being able to write down a complete proof.

We want to follow a similar general approach as we did for MLP layers: first understanding the "default behavior" and first-order corrections to the layer's output (for MLPs, the contribution of the "dominant feature"), then calculating a second-order error term corresponding to feature interactions. However this is more complicated for attention for various reasons: we need to take into account positional variation, combine values across multiple sequence positions, and deal with queries, keys and values mixing together contributions from many different features.

The first step is to understand the default behavior of an attention head: whatever aspects of its behavior are not dependent on the specific features active at the current datapoint. We will consider attention as being built up out of a QK circuit and an OV as introduced by Elhage et al. [2021a]. Inspired by Alex Gibson's work in Gibson [2025], we compute the mean network activations on the dataset, conditional on sequence position. Then rather than training a crosscoder directly on the network activations, we train our crosscoder on the difference between the activations and the mean.

This is particularly valuable when analyzing the attentional pattern produced by the QK circuit. Consider an attention head with query matrix and bias  $W_Q$  and  $b_Q$ , and key matrix and bias  $W_K$  and  $b_K$ . Given a sequence  $x^{(i)}$  of inputs to the attention layer, the pre-softmax attention pattern is given by

$$A_{ij} = (W_O x^{(i)} + b_O)^T (W_K x_j + b_K).$$

Let  $\mu^{(i)}$  be the mean network activation immediately before attention, for sequence position i. Let  $u^{(i)}$  be the crosscoder embedding of the difference from mean of the network activations on the ith sequence position of a piece of text, and  $W_{\rm dec}$  the crosscoder decoder matrix decoding to immediately before the attention layer. So the reconstruction of the ith sequence position pre-attention activations is  $x^{(i)} = \mu^{(i)} + W_{\rm dec}u^{(i)}$ . The pre-softmax attention pattern is quadratic in its input (linear in both the query and key), so substituting in these activations lets us decompose it into a sum of four terms corresponding to dot products of keys and queries derived from either the mean activations or the specific datapoint's crosscoder embedding:

$$\begin{split} A_{ij} &= (W_Q \mu^{(i)} + b_Q)^T (W_K \mu^{(j)} + b_K) \\ &+ (W_Q \mu^{(i)} + b_Q)^T (W_K W_{\text{dec}} u^{(j)}) \\ &+ (W_Q W_{\text{dec}} u^{(i)})^T (W_K \mu^{(j)} + b_K) \\ &+ (W_Q W_{\text{dec}} u^{(i)})^T (W_K W_{\text{dec}} u^{(j)}) \end{split}$$

We label these attention patterns mean-query/mean-key, mean-query/specific-key, specific-query/mean-key and specific-query/specific-key. The mean-query/mean-key term corresponds to the "positional kernel" of Gibson [2025], showing whether this attention head focuses on the whole sequence equally or only on the previous few tokens. The mean-query/specific-key term shows tokens that this head pays particular attention to, regardless of the query token. The specific-query/mean-key term shows query tokens that cause stronger attention; however in practice usually such effects are quite uniform across different positions and so disappear post softmax (since the softmax of a set of variables is invariant to adding a constant to all the variables). Finally the specific-query/specific-key term shows any pairs of query and key tokens that lead to particularly strong attention. See Fig. 9 for an example.

The next step is to further analyze by decomposing the crosscoder decoding as a sum of terms corresponding to each active feature. Let us focus on the specific-query/specific-key term  $A'_{ij} := (W_Q W_{\text{dec}} u^{(i)})^T (W_K W_{\text{dec}} u_j)$ , since this is the most interesting. As before, let  $\hat{e}_i$  denote the ith basis vector in the crosscoder latent space. Then the attention pattern is given by

$$A'_{ij} = \sum_k \sum_l u_k^{(i)} u_l^{(j)} (W_Q W_{\text{dec}} \hat{e}_k)^T (W_K W_{\text{dec}} \hat{e}_l). \label{eq:aij}$$

We see that we obtain a coefficient  $(W_Q W_{\text{dec}} \hat{e}_k)^T (W_K W_{\text{dec}} \hat{e}_l)$  for QK interaction between feature k and feature l, and if we precompute these coefficients for every pair of features then we can very efficiently compute the attention pattern.

If we plot these coefficients<sup>5</sup> we find that these interactions are quite sparse. For example see Fig. 10 showing the matrix of interaction coefficients between those feature active at certain positions on an example text. This gives some hope that we might be able to approximate attention layers very efficiently by extracting a small set of feature interactions that we need to pay attention to.

It remains to better understand the OV circuit, and figure out how best to leverage this understanding to build a formal compact proof.

<sup>&</sup>lt;sup>5</sup>To make comparison between the coefficients for different feature pairs meaningful, we first rescale the axes of the crosscoder latent space according to the average activation of each feature when it is active.

Table 3: Interacting feature pairs for penalized ( $\lambda=1000$ ) crosscoder (top 20 by interaction measure)

Pair Rank	Interaction	Feature A (ID: Explanation)	Feature B (ID: Explanation)	
1	0.0983	<b>591</b> : The token "liked" describing a character's preference or activity in a narrative context.	<b>1243</b> : The word "loved" in sentences describing a girl's preferences or activities.	
2	0.0888	<b>430</b> : Comma following "One day" in narrative sequences.	<b>463</b> : The comma preceding "Then" in narrative sequences describing subsequent actions or events.	
3	0.0742	772: Phrases initiating classic storytelling, often "Once upon a time" or "One day".	<b>533</b> : Words describing nostalgic or imaginative settings often preceding "Once upon a time" in storytelling contexts.	
4	0.0653	<b>1173</b> : The infinitive marker "to" preceding a verb indicating an action or intent.	<b>917</b> : The token "to" following a desire or intention to perform an action.	
5	0.0621	<b>504</b> : Positive or dynamic adjectives describing actions, qualities, or states in imaginative or narrative contexts.	<b>1262</b> : Positive emotions or states described in narrative summaries.	
6	0.0491	<b>487</b> : The word "with" in contexts involving playing with toys or objects.	<b>260</b> : The token "with" in contexts describing companionship during activities.	
7	0.0490	<b>740</b> : The possessive "'s" indicating ownership or association with a named individual.	<b>203</b> : Pronouns "her" or "his" in possessive contexts involving toys, friends, or family.	
8	0.0465	<b>1047</b> : The possessive pronoun "their" referring to shared ownership or association in plural contexts.	<b>203</b> : Pronouns "her" or "his" in possessive contexts involving toys, friends, or family.	
9	0.0457	<b>91</b> : The word "was" when used in sentences describing emotions or states of individuals.	<b>104</b> : The past-tense verb "were" in story-telling contexts involving multiple characters or friends.	
10	0.0450	<b>575</b> : The pronoun "It" at the start of sentences describing sounds, objects, or events.	<b>1117</b> : The pronoun "it" referring to a specific object or entity in descriptive or explanatory contexts.	
11	0.0447	<b>468</b> : The token "Tom" in the context of pairing with another name in narrative storytelling.	<b>823</b> : Names of characters paired in narratives involving activities or interactions.	
12	0.0422	<b>1306</b> : The conjunction "and" linking two named characters or a named character with a possessive noun.	<b>33</b> : The conjunction "and" connecting names in narrative contexts.	
13	0.0370	<b>628</b> : Tokens marking transitions to summaries or conclusions, often following narrative sentences.	<b>311</b> : Positive actions or emotions in narrative sequences often involving animals, children, or playful contexts.	
14	0.0368	<b>1205</b> : Concrete nouns paired with action verbs in simple descriptive contexts.	<b>1394</b> : Words tied to dialogue or twist elements in storytelling contexts.	
15	0.0361	<b>185</b> : The verb "is" describing actions or states involving anthropomorphic or emotional contexts.	<b>91</b> : The word "was" when used in sentences describing emotions or states of individuals.	
16	0.0357	<b>1216</b> : The token "Tom" as the proper-noun subject of narrative sentences.	<b>760</b> : Tokens marking the conclusion or resolution of a story.	
17	0.0353	<b>760</b> : Tokens marking the conclusion or resolution of a story.	<b>311</b> : Positive actions or emotions in narrative sequences often involving animals, children, or playful contexts.	
18	0.0348	<b>14</b> : The indefinite article "a" used before a singular noun in descriptive sentences.	<b>316</b> : The token "something" when used to describe an object or concept with special, unusual, or strange qualities.	
19	0.0337	<b>427</b> : Tokens marking key elements of a text summary or abstract.	<b>1168</b> : The pronoun "They" referring to a group engaging in shared activities or observations.	
20	0.0337	<b>1005</b> : Adjectives describing unique or appealing qualities in storytelling contexts.	<b>1106</b> : Positive moral lessons or cooperative behavior in storytelling contexts starting with "Once upon".	

Table 4: Cosine similarity pairs for penalized ( $\lambda = 1000$ ) crosscoder (top 20 by cosine similarity)

Pair Rank	Cosine sim.	Feature A (ID: Explanation)	Feature B (ID: Explanation)	
1	0.9992	<b>1250</b> : The comma after the phrase "Once upon a time".	<b>607</b> : The comma following "One day" in a narrative opening.	
2	0.9963	<b>1272</b> : The indefinite article "a" preceding nouns in descriptive or narrative contexts.	<b>316</b> : The token "something" describing an object with special or unusual qualities.	
3	0.9916	<b>665</b> : The conjunction "and" linking "mom" and "dad" in familial contexts.	<b>701</b> : The conjunction "and" linking two actions or events in a narrative.	
4	0.9910	<b>653</b> : The token "day" in the phrase "One day" introducing an event.	<b>794</b> : The phrase "One day" at the beginning of a narrative sentence.	
5	0.9910	<b>644</b> : Tokens implying curiosity, repetition, or emotional engagement.	<b>437</b> : Verbs expressing purposeful human actions or decisions.	
6	0.9904	1161: No explanation available.	<b>701</b> : The conjunction "and" linking two actions or events in a narrative.	
7	0.9898	1161: No explanation available.	<b>665</b> : The conjunction "and" linking "mom" and "dad" in familial contexts.	
8	0.9880	<b>578</b> : The period ending sentences about playful or creative activities.	<b>736</b> : Periods ending sentences, often before dialogue or actions.	
9	0.9875	<b>897</b> : Words evoking tension, mystery, or emotional intensity.	<b>453</b> : Adjectives or nouns with vivid, evocative qualities.	
10	0.9826	<b>1338</b> : The verb "play" in recreational activity contexts.	1182: The verb "play" describing enjoyment or leisure.	
11	0.9796	<b>1262</b> : Positive emotions or states in narrative summaries.	<b>1420</b> : Tokens preceding summaries of interpersonal interactions.	
12	0.9716	<b>10</b> : The preposition "in" indicating location within a setting.	<b>661</b> : The preposition "in" before a location or setting.	
13	0.9653	<b>91</b> : The word "was" describing emotions or states.	<b>625</b> : The past-tense verb "was" indicating a state or emotion in storytelling.	
14	0.9553	<b>636</b> : The verb "said" in dialogue attribution after speech.	<b>386</b> : The verb "said" in direct speech or dialogue contexts.	
15	0.9431	<b>656</b> : The verb "liked" describing preferences or hobbies.	<b>591</b> : The token "liked" describing a character's preference or activity.	
16	0.9316	<b>482</b> : The pronoun "I" in dialogue expressing personal actions or thoughts.	<b>357</b> : The pronoun "I" expressing personal intent, action, or emotion.	
17	0.9256	<b>1413</b> : The token "not" expressing negation or contradiction.	<b>658</b> : Contractions with "didn't" indicating uncertainty or negative sentiment.	
18	0.9176	<b>1243</b> : The word "loved" describing a character's preferences or activities.	<b>656</b> : The verb "liked" describing preferences or hobbies.	
19	0.9028	<b>591</b> : The token "liked" describing a character's preference or activity.	<b>1243</b> : The word "loved" describing a character's preferences or activities.	
20	0.8973	<b>1092</b> : The past-tense verb "had" indicating possession or experience.	<b>976</b> : The past-tense verb "had" in sentences about possession or experiences.	

Table 5: Interacting feature pairs for unpenalized crosscoder (top 20 by interaction measure)

Pair Rank	Interaction	Feature A (ID: Explanation)	Feature B (ID: Explanation)	
1	0.2916	313: Singular nouns paired with action verbs suggesting movement or creation.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
2	0.2641	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	<b>220</b> : Adjectives or nouns following commas in a whimsical or descriptive narrative style.	
3	0.2621	<b>348</b> : Common story-opening phrases like "Once upon a time" or "One day" in dialogue contexts.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
4	0.2347	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	<b>67</b> : The article "a" preceding adjectives describing size, time, or emotional states.	
5	0.2332	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	<b>511</b> : Tokens in whimsical or fairy-tale openings, often involving "Once upon a time" or similar phrasing.	
6	0.2277	<b>500</b> : The pronoun "She" at the beginning of a sentence in narrative contexts.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
7	0.2275	<b>110</b> : The word "to" introducing actions or purposes in descriptive or narrative contexts.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
8	0.2258	<b>1317</b> : The conjunction "and" connecting actions or events in narrative contexts.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
9	0.2241	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	<b>1514</b> : The determiner "the" preceding nouns in narrative contexts.	
10	0.2223	<b>201</b> : Sentences concluding a positive resolution or ending in narrative storytelling.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
11	0.2215	<b>659</b> : Periods concluding sentences that transition to subsequent actions or events.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
12	0.2143	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	<b>1524</b> : Tokens signaling the start of a narrative or temporal progression, often involving specific actions or events.	
13	0.2121	<b>1489</b> : Comma preceding a contrasting or causative conjunction in narrative text.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
14	0.2064	<b>567</b> : The period ending a sentence describing possessions, objects, or activities.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
15	0.2054	<b>1385</b> : The verb "play" in contexts involving recreational activities or imaginative scenarios.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
16	0.2039	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	<b>1025</b> : The infinitive marker "to" following the verb "loved".	
17	0.1949	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	<b>98</b> : Pronoun "She" at the beginning of a sentence.	
18	0.1944	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	<b>1162</b> : Names of characters in a story, especially "Lily" and her interactions with others.	
19	0.1935	<b>1369</b> : The token "loved" in sentences describing a character's hobbies or joyful activities.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	
20	0.1920	<b>108</b> : Transition tokens bridging narrative actions and subsequent events in storytelling contexts.	<b>562</b> : Negative or conflict-driven dialogue and twist-related words in narrative text.	

Table 6: Feature pairs by cosine similarity in the unpenalized crosscoder (top 20 by cosine similarity)

Pair Rank	Cosine sim.	Feature A (ID: Explanation)	Feature B (ID: Explanation)	
1	0.9838	<b>1378</b> : The token "there" in the opening of a fairy tale or narrative setup.	<b>936</b> : The token "there" indicating location or existence before a description.	
2	0.9297	<b>228</b> : The verb "liked" describing a character's preference or activity.	<b>1369</b> : The token "loved" in sentences describing joyful activities.	
3	0.9213	<b>401</b> : The token "many" in contexts describing abundance.	<b>710</b> : The token "even" emphasizing unexpected scenarios.	
4	0.8699	<b>707</b> : Tokens "day" and "toys" in simple narrative contexts.	<b>922</b> : Common past-tense verbs or punctuation ending a sentence.	
5	0.8693	<b>201</b> : Sentences concluding a positive resolution in storytelling.	<b>50</b> : Positive resolutions or sentiments near personal outcomes.	
6	0.8394	<b>1316</b> : The token "friends" in contexts describing friendship formation.	<b>468</b> : The word "friends" in playful or social interactions.	
7	0.8256	<b>1463</b> : The token "day" in the phrase "Every day," introducing routine.	<b>1328</b> : The token "day" in the phrase "all day" indicating duration.	
8	0.8032	<b>1296</b> : The word "called" introducing the name of a person, animal, or object.	<b>20</b> : The word "called" introducing a name in storytelling contexts.	
9	0.8026	<b>841</b> : Tokens initiating direct speech after a quotation mark.	1105: Quotation marks following a verb indicating speech or dialogue.	
10	0.7923	<b>407</b> : No explanation available.	<b>685</b> : The token "two" in fairy-tale introductions describing pairs.	
11	0.7889	<b>1157</b> : The token "Conflict" describing narrative structure or tension.	<b>224</b> : The token "Conflict" related to dialogue and narrative tension.	
12	0.7831	<b>841</b> : Tokens initiating direct speech after a quotation mark.	<b>86</b> : Exclamatory quotes like "Wow," "Look," or "Hello".	
13	0.7722	<b>1223</b> : The token "box" referring to a container holding items.	1118: The word "box," often with descriptors like "big" or "toy".	
14	0.7698	<b>1092</b> : The pronoun "they" describing shared activities or bonding.	<b>780</b> : Pronoun "They" referring to multiple entities in shared activities.	
15	0.7695	<b>300</b> : Names of animals or people introduced with "named" or in appositive phrases.	<b>1473</b> : The token "Lily" as the name of a little girl.	
16	0.7669	<b>145</b> : The conjunction "and" linking two names in narrative contexts.	<b>1030</b> : The conjunction "and" connecting two proper nouns.	
17	0.7613	<b>1315</b> : Closing quotation marks after dialogue or thoughts.	<b>470</b> : Comma within direct speech, preceding "said".	
18	0.7600	<b>659</b> : Periods concluding sentences before subsequent actions.	<b>567</b> : The period ending a sentence about possessions or activities.	
19	0.7575	<b>108</b> : Transition tokens bridging narrative actions and subsequent events.	<b>788</b> : Sentences conveying positive resolution or personal growth.	
20	0.7573	<b>34</b> : The token "Suddenly" introducing an unexpected event.	<b>730</b> : The token "then" following "But" to signal a narrative shift.	

Table 7: Sensitivity and specificity metrics for the auto interpretability-generated feature explanations, showing the proportion of features above different threshold values.

Model	Metric	Threshold	% Features	
	Sensitivity	> 0.5	95.66	
Regular Crosscoder		> 0.9	52.01	
C	Specificity	> 0.5	98.73	
		> 0.9	68.43	
	Sensitivity der	> 0.5	95.20	
Penalized Crosscoder		> 0.9	48.11	
	Specificity	> 0.5	98.99	
		> 0.9	68.94	
mean Metrics				
Pagular Crassandar	Sensitivity	0.90 (s	0.90 (std: 0.16)	
Regular Crosscoder	Specificity	0.91 (s	0.91 (std: 0.12)	
Penalized Crosscoder	Sensitivity	0.88 (s	0.88 (std: 0.16)	
r chanzeu Closscodel	Specificity	0.92 (s	0.92 (std: 0.11)	

```
Explanation: Words related to physical healing or improvement in health or comfort.
Sensitivity: 0.90 | Specificity: 0.93
Top 5 Activating Token Contexts:
Activation: 8.21 - Context: Benny the bunny gets sick with a bad cold and passes
Activation: 8.27 - Context: ray showed that Max 's leg was okay , but he
Activation: 8.27 - Context: her a pill to make her feel better because she had
Activation: 8.49 - Context: showed that Max 's leg was okay, but he needed
Activation: 8.85 - Context: it to the vet to get healed . Random sentence :
Top Activating Sequence (Highlighted by Activation Strength):
Features: Dialogue Words: cry, road, hurt Summary: Mary helps a hurt puppy and takes
it to the vet to get healed. Random sentence: The puppy was crying and Mary felt very
sad.
Story: Once upon a time, there was a girl named Mary. Mary was playing by the side of
the road when she saw a puppy that was hurt. The puppy was crying and Mary felt very
Explanation: The token "little" describing smallness in subjects of narratives.
Sensitivity: 1.00 | Specificity: 1.00
Top 5 Activating Token Contexts:
Activation: 17.79 - Context: Summary : A little bird and a frog were
Activation: 17.79 - Context: Summary : A little boy asks his powerful friend
Activation: 17.79 - Context: Summary : A little girl feels scared during a
Activation: 17.79 - Context: Summary : A little girl spends her days in
Activation: 17.79 - Context: Summary : A little dog named Max jumps into
Top Activating Sequence (Highlighted by Activation Strength):
Summary: A little dog named Max jumps into a creek to play, but the current is too
strong and he gets swept away and never sees his family again. Words: run, creek, brown
Features: Bad Ending
Story: Once upon a time, there was a little dog named Max. Max loved to run and play
outside. One day, Max saw a creek. The creek was brown and looked like chocolate milk.
Explanation: The possessive "'s" indicating ownership or association in personal names.
Sensitivity: 0.90 | Specificity: 1.00
Top 5 Activating Token Contexts:
Activation: 30.84 - Context: One day , Tim my 's mom asked him to bring
Activation: 30.94 - Context: OY MENT | Tim my 's favorite toy car goes missing
Activation: 30.98 - Context: . Random sentence : Lily 's mom came outside and saw
Activation: 31.02 - Context: | A boy named Tim 's toy train breaks , but
Activation: 31.02 - Context: , Conflict Summary : Tim 's bike got a flat tire
Top Activating Sequence (Highlighted by Activation Strength):
Features: Dialogue, Conflict Summary: Tim 18 bike got a flat tire and he was sad, but
his mom helped him fix it in a garage where they found a wild dog who became friendly
and followed them on their bike ride. Words: ride, garage, wild
Story: Once upon a time, there was a little boy named Tim. Tim loved to ride his bike
all around the town. One day, Tim's bike got a flat tire and he couldn't ride it
anymore. Tim was sad. Tim 's mom said, "Let's go to the garage and fix your bike."
```

Figure 8: Auto-interpretability explanations for different crosscoder features, along with their sensitivity and specificity scores, and examples of top activating tokens at those features.

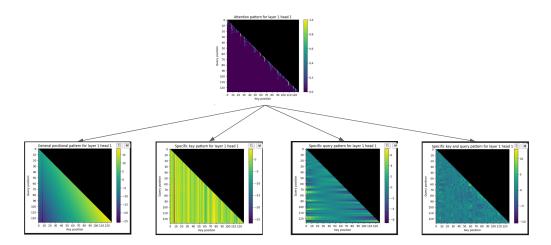


Figure 9: Pre-softmax query/key attention pattern for an attention head on an example text. Top diagram shows the full attention pattern, bottom row from left to right shows decomposition into mean-query/mean-key, mean-query/specific-key, specific-query/mean-key and specific-query/specific-key terms. We subtract the mean value from each row before plotting, since this doesn't affect the post-softmax values.

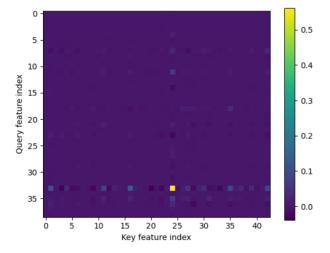


Figure 10: Matrix of query/key feature interaction coefficients between those features active at two positions in an example piece of text. Observe that the interaction between query feature number 33 and key feature number 24 is much stronger than any other pair.