

AUTOMATICALLY INTERPRETING MILLIONS OF FEATURES IN LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

While the activations of neurons in deep neural networks usually do not have a simple human-understandable interpretation, sparse autoencoders (SAEs) can be used to transform these activations into a higher-dimensional latent space which can be more easily interpretable. However, SAEs can have millions of distinct latents, making it infeasible for humans to manually interpret each one. In this work, we build an open-source automated pipeline to generate and evaluate natural language interpretations for SAE latents using LLMs. We test our framework on SAEs of varying sizes, activation functions, and losses, trained on two different open-weight LLMs. We introduce five new techniques to score the quality of interpretations that are cheaper to run than the previous state of the art. One of these techniques, intervention scoring, evaluates the interpretability of the effects of intervening on a latent, which we find explains latents that are not recalled by existing methods. We propose guidelines for generating better interpretations that remain valid for a broader set of activating contexts, and discuss pitfalls with existing scoring techniques. Our code is available online.

1 INTRODUCTION

Large language models (LLMs) have reached human level performance in a broad range of domains (OpenAI, 2023), and can even be leveraged to develop agents (Wang et al., 2023) that can strategize (Bakhtin et al., 2022), cooperate and develop new ideas (Lu et al., 2024; Shaham et al., 2024). At the same time, we understand little about the internal representations driving their behavior. Early mechanistic interpretability research focused on analyzing the activation patterns of individual neurons (Olah et al., 2020; Gurnee et al., 2023; 2024). Due to the large number of neurons to be interpreted, automated approaches were proposed (Bills et al., 2023), where a second LLM is used to propose an interpretation given a set of neuron activations and the text snippets it activates on, in a process similar to that of generating a human label. But research has found that most neurons are “polysemantic”, activating in contexts that can be significantly different (Arora et al., 2018; Elhage et al., 2022). The Linear Representation Hypothesis (Park et al., 2023) posits that human-interpretable concepts are encoded in *linear combinations* of neurons. A significant branch of current interpretability work focuses on extracting these latents and disentangling them (Bereska & Gavves, 2024).

Sparse autoencoders (SAEs) were proposed as a way to address polysemanticity (Cunningham et al., 2023). SAEs consist of two parts: an encoder that transforms activation vectors into a sparse, higher-dimensional latent space, and a decoder that projects the latents back into the original space. Both parts are trained jointly to minimize reconstruction error. SAE latents were found to be interpretable and potentially more monosemantic than neurons (Bricken et al., 2023; Cunningham et al., 2023). Recently, a significant effort was made to scale SAE training to larger models, like GPT-4 (Gao et al., 2024) and Claude 3 Sonnet (Templeton et al., 2024), and they have become an important interpretability tool for LLMs.

Training an SAE yields numerous sparse features, each of which needs a natural language interpretation. In this work, we introduce an automated framework that uses LLMs to generate an interpretation for each latent in an SAE. We use this framework to explain millions of latents across multiple models, layers, and SAE architectures. We also propose new ways to evaluate the quality of interpretations, and discuss the problems with existing approaches. We hope that these result can

054 inform other repositories of interpretations for SAE latents that already exist (Lin & Bloom, 2023),
055 as well as improve comparisons between different SAE architectures.
056

057 2 RELATED WORK 058

059 One of the most successful approaches to automated interpretability focused on explaining neurons
060 of GPT-2 using GPT-4 (Bills et al., 2023). GPT-4 was shown examples of contexts where a given
061 neuron was active and was tasked to provide a short interpretation that could capture the activation
062 patterns. To evaluate if a given interpretation captured the behavior of the neuron, GPT-4 was
063 tasked to predict the activations of the neuron in a given context having access to that interpretation.
064 The interpretation is then scored by how much the simulated activations correlate with the true
065 activations. A similar approach was used in Templeton et al. (2024) to explain SAE latents of
066 Claude 3 Sonnet. In general, current approaches focus on collecting contexts together with latent
067 activations from the model to be explained, and use a larger model to find patterns in activating
068 contexts.

069 Following those works, other methods of evaluating interpretations have been proposed, including
070 asking a model to generate activating examples and measuring how much a "module" activates
071 (Singh et al., 2023; Kopf et al., 2024). More recently, "interpretability agents" have been built,
072 managing to do iterative experiment to find the best interpretations of vision neurons (Shaham et al.,
073 2024). Potentially cheaper versions of automated interpretability have been suggested, where the
074 model that is being explained doubles as an interpretation generation model (Kharlapenko et al.,
075 2024). A prompt querying the meaning of a single placeholder token is passed to the model and
076 latent activations are patched into its residual stream at the position of the placeholder token during
077 execution, generating continuations related to the latent. This technique is inspired by earlier work
078 on Patchscopes (Ghandeharioun et al., 2024) and SelfIE (Chen et al., 2024).
079

080 3 AUTOMATED INTERPRETABILITY PIPELINE 081

082 In this Section we explain, step-by-step, the main pipeline used to generate interpretations and score
083 them. The pipeline can be broadly divided into 3 sequential steps. First, the activations of the SAEs
084 to be interpreted are collected over a broad dataset. Then, for each latent, relevant contexts are se-
085 lected and shown to an LLM which generates an interpretation for the activation pattern observed.
086 Finally, these interpretations are matched with different contexts and used by an LLM in different
087 tasks that evaluate how good the interpretations are in predicting activating and non-activating con-
088 texts, details explained below. As an illustration, we represent how that pipeline might look like for
089 a real latent, see Figure 1.

090 3.1 COLLECTING ACTIVATIONS 091

092 Sparse autoencoders (SAEs) are feedforward networks with one hidden layer trained to reproduce
093 their input using a sparse number of neurons. This means that while there are several times more
094 neurons than there are dimensions of the input, only a small fraction of these neurons— less than
095 60 in this work— are non-zero and contribute to the output. In our pipeline, we seek to explain the
096 non-zero activations of the SAE hidden layer neurons, which we call "latents."

097 We collected latent activations from the SAEs over a 10M token sample of RedPajama-v2 (RPJv2;
098 Computer 2023). This is the same dataset that we used to train the Llama 3.1 8b SAEs, and consists
099 of a data mix similar to the Llama 1 pretraining corpus.

100 We collected batches of 256 tokens starting with the beginning of sentence token (BOS).¹ The con-
101 texts used for activation collection are smaller than the contexts used to train the SAEs, and we find
102 that on average, 30% of the latents of the 131k latent, per layer, Gemma 2 9b SAE don't activate
103 more than 200 times over these 10M tokens and 15% don't activate at all. When we consider the
104 training context length of 1024, only 5% of latents don't fire. When using a closer proxy of the
105 training data, the "un-copyrighted" Pile, we find that the number of latents that activate fewer than
106

107 ¹We throw out the activations on the BOS token when generating interpretations for Gemma, as we were
told in personal communication that these SAEs were not trained on BOS activations.

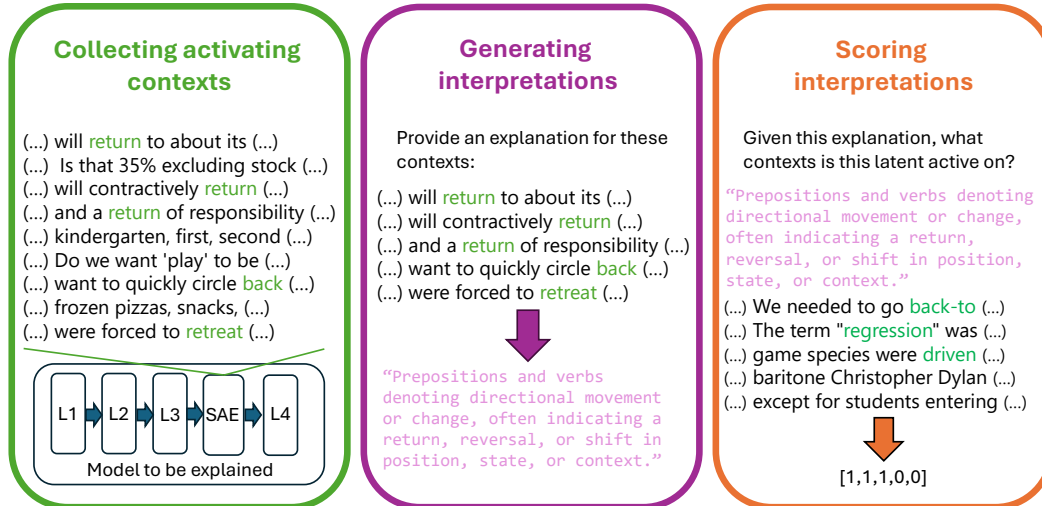


Figure 1: **Auto-interpretability pipeline.** The first step of the interpretability pipeline is collecting the activations of the target SAE over a broad range of text. In this figure we show short contexts for legibility, and over represent the active latent, which in reality would be active in a very small portion of the text. The activating contexts are selected and shown to an explainer LLM, prompt in Appendix A.2, which provides a short interpretation. This interpretation is then given to a scorer LLM that is tasked to use this interpretation to distinguish activating from non activating contexts, see Section 3.3

200 times decreases to 15%, and that only 1% don’t activate at all, even when considering contexts of size 256. Interestingly, in the case of the 16k latent Gemma 2 9b SAE only around 10% of latents activate fewer than 200 times on the RPJv2 sample, suggesting that the larger 131k latent SAE learned more dataset specific latents than its smaller cousin.

3.2 GENERATING INTERPRETATIONS

Our approach to generating interpretations follows Bricken et al. (2023) in showing the explainer model examples sampled from different quantiles, but uses a more “natural” prompt where the activating example is shown whole, with the activating tokens emphasized and their activation strength shown after the example, see Appendix A.2 for the full prompt. For examples of some interpretations and their scores, see Figure A1.

We show the explainer model, Llama 3.1 70b Instruct, 40 different activating examples, each activating example consisting of 32 tokens, where the activating tokens can be on any position. While most latents are well-explained by contexts of length 32, we expect latents active over longer contexts to not be well captured by our pipeline. Using longer contexts to generate explanations would require either showing the model less examples, or using a much more powerful (and expensive) explainer model. Because we are able to evaluate interpretations significantly cheaper, bad interpretation can easily be filtered out, and we leave the explanation of “long-range” latents for future work. Using this interpretation setup would cost \$200 to interpret 1 million latents.

We find that randomly sampling from a broader set of examples leads to interpretations that cover a more diverse set of activating examples, sometimes to the detriment of the top activating examples, see Figure 3. Using only top activating examples often yields more concise and specific interpretations that accurately describe these examples, but which fail to capture the whole distribution. On the other hand, stratified sampling from the deciles of the activation distribution can lead to interpretations that are too broad and that fail to capture a meaningful interpretation. For examples of these types of failure modes, see Figure A2 and the discussion in the Appendix A.3.

3.2.1 AUTOMATICALLY INTERPRETING INTERVENTIONS

Automatic interpretability methods, including most of those explored in this paper, typically look for correlations between activation of a latent and some natural-language property of the input. However, some latents are more closely related to what the model will output. For example, we find a latent² whose activation causes the model to output words associated with reputation but does not have a simple interpretation in terms of inputs.

We define an *output latent* as a latent that causes a property of the model’s output that can be easily explained in natural language. See Section 3.3.5 for the definition we use in scoring.

Output latents can also be described in terms of their correlation with inputs. For example, the “reputation” latent activates in contexts where likely next tokens relate to reputation. However, explaining output latents in terms of causal influence on output has two advantages.

1. **Scalability.** Output latents are easier to describe in terms of effect on output because the explainer only needs to notice a simple pattern of output. The pattern of inputs this latent correlates with is more complex. Explaining a latent by correlating it with inputs requires approximating the computation of the subject model leading up to that latent, which may be challenging when subject models are highly capable and performing difficult tasks. Some latents’ influence on output, however, might remain easily explainable.
2. **Causal evidence.** We might like to know that a latent causes a property of the model’s output so that we can steer the model. Previous work has shown that existing auto-interpretability scores fail to accurately capture how well a given interpretation can predict the effect of intervening on a given neuron (Huang et al., 2023). Further, prior work has argued that causal evidence is more robust to distribution shifts (Bühlmann, 2018; Schölkopf et al., 2012).

3.3 SCORING INTERPRETATIONS

Not only should the generation of interpretations efficient, but so should the scoring of such interpretations. Practitioners would like to know how faithful each interpretation is to the actual behavior of the network. Some latents may not be amenable to a simple interpretation, and in these cases we expect the auto-interpretability pipeline to output an interpretation with poor evaluation metrics. Secondly, we can use evaluations as a feedback signal for the training of explainer models and tuning hyperparameters in the pipeline. Finally, some SAEs may simply be poor quality overall, and the aggregate evaluation metrics of the SAE’s interpretations can be used to detect this.

The quality of SAE interpretations has so far been measured via simulation scoring, a technique introduced by Bills et al. (2023) for evaluating interpretations of neurons. It involves asking an explainer model to “simulate” a latent in a set of contexts, which means predicting how strongly the latent should activate at each token in a context given an interpretation. The Pearson correlation between the simulated and true activations is the simulation score. The standard in the literature is to sample contexts from a “top-and-random” distribution that mixes maximally activating contexts and contexts sampled uniformly at random from the activating corpus. While oversampling the top activating contexts introduces bias, it is used as a cheap variance reduction technique, given that simulation scoring over hundreds of examples per latent would be expensive (Cunningham et al., 2023).

In this work, we take a slightly different view of what makes for a good interpretation: an interpretation should serve as a *binary classifier* distinguishing activating from non-activating contexts. The reasoning behind this is simple. Given the highly sparse nature of SAE latents, most of their variance could be captured by a binary predictor that predicts the mean nonzero activation when the latent is expected to be active, and zero otherwise. In statistics, zero-inflated phenomena are often modeled as a mixture distribution with two components: a Dirac delta on zero, and another distribution (e.g. Poisson) for nonzero values. This goes against the general wisdom of simulation scoring, which focuses only on activating examples, even though they are less than 0.01% of the relevant contexts

²latent 157 of Gemma 2 9b’s layer 32 autoencoder with 131k latents and an average L0 norm of 51. It has a detection score of 0.6.

for each latent of an SAE and that being able to distinguish non-activating contexts from activating contexts seems to be relatively more important.

Detection	Fuzzing
<p>Explanation: Words related to American football positions, specifically the tight end position</p> <p>Sentences: 1: Patriots tight end Rob Gronkowski had his boss 2: names of months used in The Lord of the Rings 3: shown, is generally not eligible for ads. For example</p> <p>Correct output: [1,0,0]</p>	<p>Explanation: Words related to American football positions, specifically the tight end position</p> <p>Sentences: 1: Patriots <tight end> Rob Gronkowski had his boss 2: You should know this <about> offensive line coaches 3: <running backs>," he said. .. Defensive<end></p> <p>Correct output: [1,0,1]</p>
Surprisal	Embedding
<p>Explanation 1: Words related to American football positions, specifically the tight end position.</p> <p>Explanation 2: Sentences about dogs</p> <p>Sentence: Patriots tight end Rob Gronkowski had his boss -</p> <p>Correct output: $P(\text{Sentence} \text{Exp 1}) > P(\text{Sentence} \text{Exp 2})$</p>	<p>Explanation: Words related to American football positions, specifically the tight end position.</p> <p>Sentence 1: names of months used in The Lord of the Rings</p> <p>Sentence 2: Patriots tight end Rob Gronkowski had his boss -</p> <p>Correct output: $\text{Embed}(S2) \cdot \text{Embed}(\text{Exp}) > \text{Embed}(S1) \cdot \text{Embed}(\text{Exp})$</p>

Figure 2: **The new proposed scoring methods.** In **detection** scoring, the scorer model is tasked with selecting the set of sentences that activate a given latent given an interpretation. In this work, we show 5 examples at the same time, and each has an identical probability of being a sentence that activates the latent, independent of whether any other example also activates the latent. The activating tokens are colored in green for display, but that information is not shown to the scorer model. For **fuzzing** scoring, the scorer model is tasked with selecting the sentences where the highlighted tokens are the tokens that activate a target latent given an interpretation of that latent. In **surprisal** scoring, activating and non-activating examples are run through the model and the loss over those sentences is computed. Correct interpretations should decrease the loss in activating sentences compared to a generic interpretation, but shouldn’t significantly decrease the loss in non-activating sentences. For **embedding** scoring, activating and non-activating sentences are embedded as “documents” that should be retrieved using the interpretation as a query.

As an alternative to simulation scoring, we introduce four new evaluation methods that focus on how well an interpretation enables a scorer to discriminate between activating and non-activating contexts. As an added benefit, all of these methods are more compute-efficient than simulation, see Table 1. We also propose a scoring method that evaluates the interpretability of interventions using a specific latent.

3.3.1 DETECTION

One simple approach for scoring interpretations is to ask a language model to identify whether a whole sequence activates a SAE latent given an interpretation. Detection requires few output tokens for each example used in scoring, meaning examples from a wider distribution of latent activation strengths can be used at approximately the same expense. By including non-activating contexts, this method measures both the precision and recall of the interpretation.

Detection is more “forgiving” than simulation insofar as the scorer does not need to localize the latent to a particular token. This means that detection evaluates whether an interpretation correctly identifies the types of contexts that a latent is active on, even if the interpretation does not correctly identify the actual token.

Both this method and fuzzing (described next) can leverage token probabilities to estimate how certain the scorer model is of their classification, and this is an effect that we believe can be to improve the scoring methods. Details on the prompt in Appendix A.4.1.

3.3.2 FUZZING

Fuzzing is similar to detection, but at the level of individual tokens. Here, potentially activating tokens are delimited in each example and the language model is prompted to identify which of the sentences are correctly marked. Fuzzing is the most similar with simulation. Because SAE activations are very sparse, simulation scoring mostly boils down to correctly identifying which of the tokens have non-zero activations. For this reason, fuzzing is the score that most correlates with simulation.

When interpretations focus on which tokens the latent activates, but not on which contexts, they can still score high on fuzzing, but comparatively lower on detection, specially if they are common tokens. Thus, evaluating an interpretation on both detection and fuzzing can identify whether a model is classifying examples for the correct reason. Details on the prompt on Appendix A.4.2.

3.3.3 SURPRISAL

Surprisal scoring is based on the idea that a good interpretation should help a base language model \mathcal{M} (the “scorer”) achieve lower cross-entropy loss on activating contexts than it would without a relevant interpretation. Specifically, for each context \mathbf{x} , activating or non-activating, we measure the *information value* of an interpretation \mathbf{z} as $\log p_{\mathcal{M}}(\mathbf{x}|\mathbf{z}) - \log p_{\mathcal{M}}(\mathbf{x}|\tilde{\mathbf{z}})$, where $\tilde{\mathbf{z}}$ is a fixed pseudo-interpretation. A good interpretation should have higher information value on activating examples than on non-activating ones.

The overall surprisal score of the interpretation is given by the AUROC of its information value when viewed as a classifier distinguishing activating from non-activating contexts. Details on the prompt and how to compute the score in Appendix A.4.3.

Surprisal relies on the ability of the scorer model to tune its predictions based on the explanation, measuring both the relevance of the context given in the interpretation and the token specificity, but we believe that the current setup could be improved, as this is the score with the least correlation with the others.

3.3.4 EMBEDDING

Classifying between active and non-active contexts given a certain interpretation can also be seen as using interpretations of latents as “queries” that should be able to retrieve relevant “documents”, contexts where the latent is active, between non-relevant “documents”, non-activating contexts. This way, we take a selection of activating and non-activating contexts that embedded by an encoding transformer, and the similarity between the query and the documents is used as a classifier to distinguish between activating and non-activating contexts, and the score is given by the AUROC.

If the encoding model is small enough - we used a 400M parameter model - this technique is the fastest and opens up the possibility to evaluate a larger fraction of the activation distribution. We have seen that using a larger embedding model - 7B parameter model - didn’t significantly improve the scores, see fig A3, although we believe that this approach was under-investigated. Details on the prompt, on the embedding model and on the way to compute the score in Appendix A.4.4.

3.3.5 INTERVENTION SCORING

Unlike the above four context-based scores, intervention scoring interprets a latent’s counterfactual impact on model output. We quantify the interpretability of an intervention I with interpretation \mathbf{z} on a distribution of prompts π as the average decrease in the scorer’s surprisal about the interpretation when conditioned on text generated with the intervention.

$$S(I, \mathbf{z}; \pi) = \mathbb{E}_{\mathbf{x} \sim \pi} \left[\mathbb{E}_{\mathbf{i} \sim \mathcal{G}_I(\mathbf{x})} [\log p_{\mathcal{M}}(\mathbf{z}|\mathbf{i})] - \mathbb{E}_{\mathbf{g} \sim \mathcal{G}(\mathbf{x})} [\log p_{\mathcal{M}}(\mathbf{z}|\mathbf{g})] \right] \quad (1)$$

$\mathcal{G}(\mathbf{x})$ is the distribution over subject model generations given prompt \mathbf{x} at temperature 1. In $\mathcal{G}_I(\mathbf{x})$, the intervention is applied to the subject model as it generates. In practice we estimate the quantity S by sampling one clean and one intervened generation for each sampled of prompt.

Sufficiently strong interventions can be trivially interpretable by causing the model to deterministically output some logit distribution. Therefore, **interpretability scores of interventions should be**

324 **compared for interventions of a fixed strength.** We define the strength σ of an arbitrary interven-
 325 tion I as the average KL-divergence of the model’s intervened logit distribution with reference to
 326 the model’s clean output.

$$327 \sigma(I; \pi) = \mathbb{E}_{\mathbf{x} \sim \pi} [D_{KL}(p_{\text{subject}}(\cdot | \mathbf{x}) || p_{\text{subject}, I}(\cdot | \mathbf{x}))] \quad (2)$$

329 See Appendix A.7 for details on the interpretation and scoring pipeline we use in our intervention
 330 experiments.

333 4 RESULTS

335 4.1 COMPARING SCORING METHODS

336
 337 When wanting to evaluate the explanation over a larger number of examples, simpler methods than
 338 simulation scoring have to be used (Templeton et al., 2024). Simulation scoring is traditionally
 339 done only in activating examples, ignoring whether the proposed explanation correctly handles non-
 340 activating examples. Simulation scoring is normally done in a small amount of examples per latent
 341 due to its high cost, and this is a significant disadvantage that prompted us to investigate more
 342 efficient scoring methods.

343 In Table 1 we compare the amount of tokens used to score a single latent using 100 different con-
 344 texts, for both fuzzing, detection and two different simulation methods. We can see that fuzzing
 345 and detection are at least 5x cheaper than simulation when using the all-at-once (AAO) method
 346 described in Bills et al. (2023), although this method requires access to the log-probabilities of the
 347 prompt tokens, which are not accessible from providers of state-of-the-art closed-source models. A
 348 researcher wanting to generate simulation scores on the fly, without access to a local model, would
 349 have to use the token-by-token (TBT) method, which is 30x more than expensive than both fuzzing
 350 and detection. For the same set of interpretations, fuzzing and detection scores given by Llama 70b
 351 and Claude Sonnet 3.5 have similar distributions. On the other hand, Claude Sonnet 3.5 produces to
 352 higher simulation scores on average than Llama 70b, see Table A8 meaning that we should probably
 353 compare the price of fuzzing and detection on Llama 70b to the price of simulation using Claude
 354 Sonnet 3.5.

355 Embedding scoring is even cheaper, requiring 4000 input tokens for 100 different contexts, which at
 356 an average of \$0.13 per million tokens would only cost \$50/100k latents.

357
 358 Table 1: Estimated cost of the different scoring methods. Brackets represent tokens that can be
 359 cached, potentially saving costs, although we do not do these calculations here. The number of
 360 input and output tokens is computed considering 100 examples show, and using fewer examples
 361 would cost proportionally less.

	Input tokens	Output tokens	Cost (\$/100k latents) Llama 70b	Cost (\$/100k latents) Claude Sonnet 3.5
Fuzzing	19.6k (14.2k)	249	676	6.2k
Detection	17.0k (11.9k)	240	588	5.5k
Simulation (AAO)	104.9k (87.5k)	5	3.6k	31.5k
Simulation (TBT)	496.9k (451k)	46.7k	18.7k	219.1k

362
 363
 364 We propose that latent interpretations should be evaluated using more than a single technique when
 365 possible, because as discussed in Section 3.3, each of the proposed methods have different failure
 366 modes, and a more rich scoring framework can address some of these flaws. Specifically, we find
 367 that when computing the correlation of scores computed on the same 800 explanations, we find that
 368 while there is a clear positive correlation, some scores disagree– see Appendix for more examples.

369 We find fuzzing to be the most correlated with detection, and this is due to fact that fuzzing is
 370 mostly measuring how well the explanation can help the model predict which tokens have non-zero
 371 activation. Due to the sparsity of SAE activations, simulation scoring is doing a very similar test. On
 372 the contrary, we find that simulation has a significantly lower correlation with detection, embedding
 373 and surprisal scoring, as they more accurately measure whether a context activates a given latent,
 374
 375
 376
 377

378 instead of which token is active in a given context. When comparing the simulation scores provided
 379 by Claude 3.5 Sonnet, all correlations have a slight increase, see A.4.6.

380
 381 Due to their relative cheapness compared with simulation scoring, we propose that fuzzing and
 382 detection scoring can be used to estimate whether explanations correctly identify on which tokens a
 383 latent is most likely to be active, and in which types of contexts this happens. Cheaper methods like
 384 embedding can be used to quickly iterate on explanations or to broadly separate latents that have
 385 bad explanations from those with good explanations, so that they can be refined.

386
 387 Table 2: Spearman correlation, computed over scores of 800 different latent interpretations. For
 388 details on how the scores are computed and for the Pearson correlation, see Appendix A.4.6

	Fuzzing	Detection	Simulation	Embedding	Surprisal
Fuzzing	1	0.73	0.75	0.41	0.30
Detection		1	0.44	0.71	0.62
Simulation			1	0.28	0.15
Embedding				1	0.79
Surprisal					1

396 397 4.1.1 INTERVENTION SCORING

399 The previously discussed scoring methods are *correlational* in the sense that they measure whether
 400 an interpretation can be used to predict the activation values of a given latent, or distinguish between
 401 activating and non-activating examples. Intervention scoring proposes to measure how well a given
 402 interpretation can predict the effect of *interventions* on the corresponding latent. Here we compare
 403 correlational interpretations generated with our pipeline and scored with fuzzing, to a set of inter-
 404 ventional interpretations scored with intervention scoring. Our hypothesis is that some latents will
 405 have low correlational scores because their behavior is better explained by their downstream effects
 406 than by the contexts where they are active.

407 In Figure 4, we see that there is a slight negative correlation between the fuzzing score and the in-
 408 tervention score, showing that on average, latents with low fuzzing score can be better explained
 409 using their downstream effects than latents with high fuzzing score. We also find that the distri-
 410 bution of intervention scores of interpretation on a trained SAE is significantly different from the
 411 distribution of an interpretation of a randomly initialized SAE or a real SAE with randomly assigned
 412 interpretations, supporting the validity of this scoring method.

413 414 4.2 COMPARING INTERPRETATION METHODS

415 We use a set of 500+ latents as a testbed to measure the effects of design choices and hyperparam-
 416 eters on interpretation quality. Each latent is scored using 100 activating and 100 non-activating
 417 examples. The activating examples are chosen via stratified sampling such that there are always
 418 10 examples from each of the 10 deciles of the activation distribution. We evaluate interpretation
 419 quality using fuzzing, detection and embedding scores, as those were both quick to compute and
 420 easy to interpret. We expect this mix of scores reflect the extent to which the proposed interpretation
 421 is valid over both activating and non-activating examples.

422 We find that the interpretations generated by considering only the activating contexts are signifi-
 423 cantly different from those generated by selecting activating contexts from the whole distribution.
 424 Interpretations based on top examples have higher specificity, but lower sensitivity. In fact, we
 425 observe that the sensitivity depends on how the degree of activation of each example shown, and
 426 that this sensitivity decreases faster for interpretation based on top examples. Interpretations gener-
 427 ated from top activating examples are more concise and specific interpretations but fail to capture
 428 the whole distribution, while explanations based on examples sampled from the whole distribu-
 429 tion lead to interpretations that are too broad, see Appendix A.3. This effect would not be seen if
 430 the scoring were done on just the most activating examples, underscoring a problem with current
 431 auto-interpretability evaluations, which produce interpretations using top activating examples and
 evaluate them on a small subset of the activation distribution.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

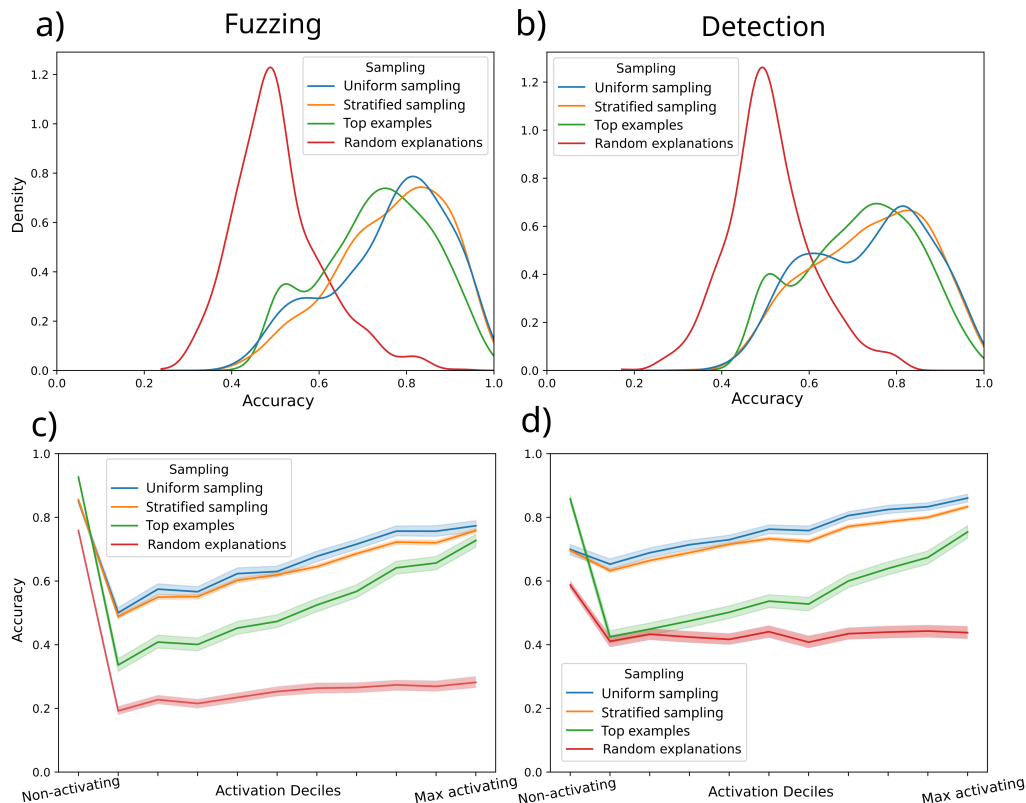


Figure 3: **Fuzzing and detection scores for different sampling techniques.** Panels a) and b) show the distributions of fuzzing and detection scores, respectively, as a function of different example sampling methods for interpretation generation. Sampling only from the top activation gets on average lower accuracy in fuzzing and on detection when compared with uniform sampling and stratified sampling. The distributions from random sampling and sampling from quantiles are very similar. Panels c) and d) measure how the interpretations generalize across activation quantiles, showing that interpretations generated from the top quantiles are better at distinguishing non-activating examples, but have lower accuracy on other quantiles, especially on the lower activating deciles. This also happens for the other interpretations, but the accuracy does not drop as much in lower activating deciles. We also show the scores of random explanations.

Increasing the size of the explainer model increases the scores of the interpretations, but we don't find the explanations generated by Claude Sonnet 3.5 to have much higher scores than those generated by Llama 3.1 70b, see A8, and both are similar to those generated by a human. Not surprisingly, we also see that even for simpler scores like fuzzing and detection, using a smaller scorer model leads to lower scores. It is possible that the explanations generated by Claude would be better than those generated by Llama 3.1 70b had we optimized the prompting techniques for that model.

Showing the explainer model a larger number of examples leads to slightly higher scores (Table A6). We find it doesn't matter much whether we use the same dataset as the training set of the SAEs (Table A2) or whether we slightly change the size of the contexts shown (Table A4). Using COT, at least when generating interpretations using Llama 70b, does not significantly increase their quality (Table A3) while significantly increasing the compute and time required to generate them. For this reason, we have not used it for our main experiments.

We find that SAEs with more latents have higher scores, and scores that are significantly higher than those of neurons. Neurons are more interpretable if made sparser by only considering the top k most activated neurons on a given token, but still significantly underperform SAEs in our tests (Table A9). The location of the SAE matters; residual stream SAEs have slightly better scores than ones trained

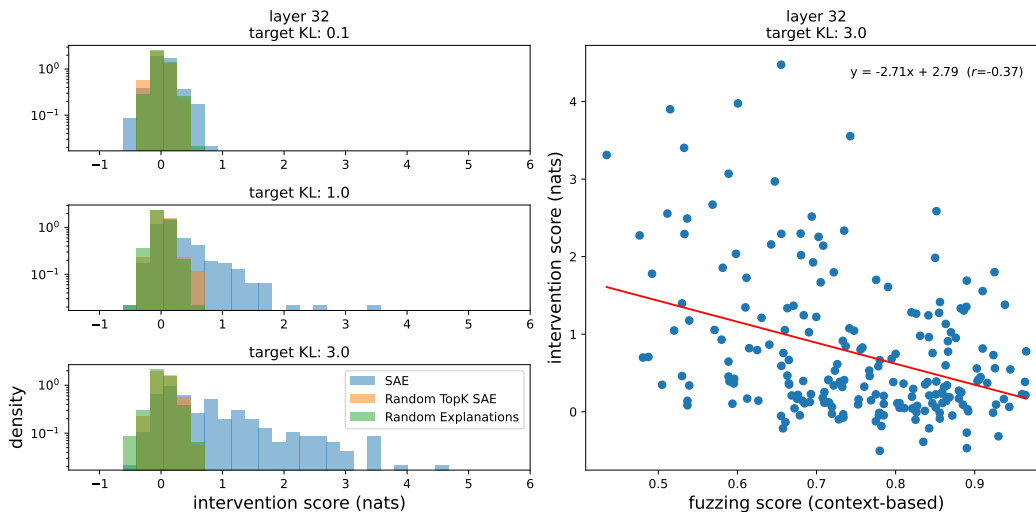


Figure 4: **Intervention scores.** Here we present intervention scores (Sec 3.3.5) for SAE latents in Gemma 2 9B at layer 32. **Left:** SAE latents are more interpretable than random latents, especially when intervening more strongly. Our explainer also produces interpretations that are scored higher than random interpretations. **Right:** Many latents that would normally be uninterpreted when using context-based automatic interpretability are interpretable in terms of their effects on output.

on MLP outputs (Table A9). We also observe that earlier layers have lower overall scores, but that the scores across model depth remain constant after the first few layers (Fig A5).

5 CONCLUSION

Explaining the latents of SAEs trained on cutting-edge LLMs is a computationally demanding task, requiring scalable methods to both generate interpretations and assess their quality. We addressed issues with the conventional simulation-based scoring and introduced five new scoring techniques, each with distinct strengths and limitations. These methods allowed us to explore the “prompt design space” for generating effective interpretations and propose practical guidelines.

Additionally, although current scoring methods do not account for interpretation length, we believe shorter interpretations are generally more useful and will incorporate this in future metrics. Some scoring methods also require further refinement, particularly in selecting non-activating examples to improve evaluation. We also investigated a method to generate and scoring interpretations of latents that is based on their downstream effect and shown that low score “correlational” scores could be due to the existence of “output” features.

Access to better, automatically generated interpretations could play a crucial role in areas like model steering, concept localization, and editing. We hope that our efficient scoring techniques will enable feedback loops to further enhance the quality of interpretations.

REFERENCES

- Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. Linear algebraic structure of word senses, with applications to polysemy. *Transactions of the Association for Computational Linguistics*, 6:483–495, 2018.
- Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, Athul Paul Jacob, Mojtaba Komeili, Karthik Konath, Minae Kwon, Adam Lerer, Mike Lewis, Alexander H. Miller, Sandra Mitts, Adithya Renduchintala, Stephen Roller, Dirk Rowe, Weiyan Shi, Joe Spisak, Alexander Wei, David J. Wu, Hugh Zhang, and Markus Zijlstra. Human-level play in the game of diplomacy by com-

- 540 binning language models with strategic reasoning. *Science*, 378:1067 – 1074, 2022. URL
541 <https://api.semanticscholar.org/CorpusID:253759631>.
542
- 543 Leonard Bereska and Efstratios Gavves. Mechanistic interpretability for ai safety—a review. *arXiv*
544 *preprint arXiv:2404.14082*, 2024.
- 545 Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever,
546 Jan Leike, Jeff Wu, and William Saunders. Language models can explain neurons in lan-
547 guage models. URL [https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.](https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html)
548 [html](https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html).(Date accessed: 14.05. 2023), 2, 2023.
549
- 550 Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Con-
551 erly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu,
552 Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex
553 Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter,
554 Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language
555 models with dictionary learning. *Transformer Circuits Thread*, 2023. [https://transformer-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)
556 [circuits.pub/2023/monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 557 Peter Bühlmann. Invariance, causality and robustness. *arXiv preprint arXiv:1812.08233*, 2018.
558
- 559 Haozhe Chen, Carl Vondrick, and Chengzhi Mao. Selfie: Self-interpretation of large language model
560 embeddings, 2024. URL <https://arxiv.org/abs/2403.10949>.
- 561 Together Computer. Redpajama: an open dataset for training large language models, 2023. URL
562 <https://github.com/togethercomputer/RedPajama-Data>.
563
- 564 Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-
565 coders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*,
566 2023.
567
- 568 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,
569 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposi-
570 tion. *arXiv preprint arXiv:2209.10652*, 2022.
- 571 Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya
572 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint*
573 *arXiv:2406.04093*, 2024.
574
- 575 Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. Patchscopes:
576 A unifying framework for inspecting hidden representations of language models, 2024. URL
577 <https://arxiv.org/abs/2401.06102>.
- 578 Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bert-
579 simas. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint*
580 *arXiv:2305.01610*, 2023.
581
- 582 Wes Gurnee, Theo Horsley, Zifan Carl Guo, Tara Rezaei Kheirkhah, Qinyi Sun, Will Hathaway,
583 Neel Nanda, and Dimitris Bertsimas. Universal neurons in gpt2 language models. *arXiv preprint*
584 *arXiv:2401.12181*, 2024.
- 585
- 586 Jing Huang, Atticus Geiger, Karel D’Oosterlinck, Zhengxuan Wu, and Christopher Potts. Rig-
587 orously assessing natural language explanations of neurons. *arXiv preprint arXiv:2309.10312*,
588 2023.
- 589 Caden Juang, Gonçalo Paulo, Jacob Drori, and Belrose Nora. Understanding and steering Llama
590 3, 9 2024. URL <https://goodfire.ai/blog/research-preview/>.
591
- 592 Dmitrii Kharlapenko, neverix, Neel Nanda, and Arthur Conmy. Self-explaining SAE fea-
593 tures, 8 2024. URL [https://www.lesswrong.com/posts/8ev6coxChSWcxCDy8/](https://www.lesswrong.com/posts/8ev6coxChSWcxCDy8/self-explaining-sae-features)
[self-explaining-sae-features](https://www.lesswrong.com/posts/8ev6coxChSWcxCDy8/self-explaining-sae-features).

- 594 Laura Kopf, Philine Lou Bommer, Anna Hedström, Sebastian Lapuschkin, Marina M. C. Höhne,
595 and Kirill Bykov. Cosy: Evaluating textual explanations of neurons, 2024. URL [https://](https://arxiv.org/abs/2405.20331)
596 arxiv.org/abs/2405.20331.
597
- 598 Johnny Lin and Joseph Bloom. Neuronpedia: Interactive reference and tooling for analyzing neural
599 networks with sparse autoencoders, 2023. URL <https://www.neuronpedia.org>. Soft-
600 ware available from neuronpedia.org.
- 601 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist:
602 Towards fully automated open-ended scientific discovery, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2408.06292)
603 [abs/2408.06292](https://arxiv.org/abs/2408.06292).
604
- 605 Tom Macgrath. Understanding and steering Llama 3, 9 2024. URL [https://goodfire.ai/](https://goodfire.ai/blog/research-preview/)
606 [blog/research-preview/](https://goodfire.ai/blog/research-preview/).
- 607 Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. Mteb: Massive text embed-
608 ding benchmark. *arXiv preprint arXiv:2210.07316*, 2022. doi: 10.48550/ARXIV.2210.07316.
609 URL <https://arxiv.org/abs/2210.07316>.
- 610 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.
611 Zoom in: An introduction to circuits. *Distill*, 5(3), March 2020. ISSN 2476-0757. doi: 10.23915/
612 [distill.00024.001](http://dx.doi.org/10.23915/distill.00024.001). URL <http://dx.doi.org/10.23915/distill.00024.001>.
613
- 614 OpenAI. Gpt-4 technical report, 2023. URL <https://arxiv.org/abs/2303.08774>.
615
- 616 Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the
617 geometry of large language models. *ArXiv*, abs/2311.03658, 2023. URL [https://api.](https://api.semanticscholar.org/CorpusID:265042984)
618 [semanticscholar.org/CorpusID:265042984](https://api.semanticscholar.org/CorpusID:265042984).
- 619 Bernhard Schölkopf, Dominik Janzing, Jonas Peters, Eleni Sgouritsa, Kun Zhang, and Joris Mooij.
620 On causal and anticausal learning. *arXiv preprint arXiv:1206.6471*, 2012.
- 621 Tamar Rott Shaham, Sarah Schwettmann, Franklin Wang, Achyuta Rajaram, Evan Hernandez, Jacob
622 Andreas, and Antonio Torralba. A multimodal automated interpretability agent, 2024. URL
623 <https://arxiv.org/abs/2404.14394>.
624
- 625 Chandan Singh, Aliyah R. Hsu, Richard Antonello, Shailee Jain, Alexander G. Huth, Bin Yu, and
626 Jianfeng Gao. Explaining black box text modules in natural language with language models,
627 2023.
- 628 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen,
629 Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L
630 Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers,
631 Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan.
632 Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Trans-*
633 *former Circuits Thread*, 2024. URL [https://transformer-circuits.pub/2024/](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html)
634 [scaling-monosemanticity/index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).
- 635 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan,
636 and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models,
637 2023. URL <https://arxiv.org/abs/2305.16291>.
638
639
640
641
642
643
644
645
646
647

A APPENDIX

A.1 SAE INTERPRETATIONS OF A RANDOM SENTENCE

In this section, we show some non cherry-picked examples of latent interpretations on a sentence.



Figure A1: **SAE latents interpretations for a random sentence.** To visualize the latent interpretations produced, we select a sentence taken from the RPJv2 dataset. We selected 4 tokens in different positions in that sentence and filter for latents that are active in different layers. Then we randomly select active latents and their corresponding interpretations to display. We display the detection and fuzzing scores of each interpretation, which indicate how well it explains other examples in the dataset (see Section 3 for details on these scores). The latents selected had high activation, but were not cherry-picked based on interpretations or scores.

A.2 EXPLAINER PROMPT

The system prompt if not using Chain of Thought (COT):

You are a meticulous AI researcher conducting an important investigation into patterns found in language. Your task is to analyze text and provide an interpretation that thoroughly encapsulates possible patterns found in it.

Guidelines:

You will be given a list of text examples on which special words are selected and between delimiters like << this >>.

If a sequence of consecutive tokens all are important, the entire sequence of tokens will be contained between delimiters <<just like this>>. How important each token is for the behavior is listed after each example in parentheses.

- Try to produce a concise final description. Simply describe the text latents that are common in the examples, and what patterns you found.
- 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.

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

- Do not mention the marker tokens ($\$ \ll \$ \gg \$$) in your interpretation.

- Do not make lists of possible interpretations. Keep your interpretations short and concise.

- The last line of your response must be the formatted interpretation, using [interpretation]:

We add to the previous prompt the following if we want to do COT:

To better find the interpretation for the language patterns, go through the following stages:

1. Find the special words that are selected in the examples and list a couple of them. Search for patterns in these words, if there are any. Don't list more than 5 words.
2. Write down general shared latents of the text examples. This could be related to the full sentence or to the words surrounding the marked words.
3. Formulate a hypothesis and write down the final interpretation using [interpretation]:.

One of the few shot examples of how examples are displayed to the model.

Example 1: and he was \ll over the moon \gg to find

Activations: (``over", 5), (`` the", 6), (`` moon", 9)

Example 2: we'll be laughing \ll till the cows come home \gg ! Pro

Activations: (``till", 5), (`` the", 5), (`` cows", 8), (`` come", 8), (`` home", 8)

Example 3: thought Scotland was boring, but really there's more \ll than meets the eye \gg ! I'd

Activations: (``than", 5), (`` meets", 7), (`` the", 6), (`` eye", 8)

If COT is used, an explicit example of using COT is demonstrated.

ACTIVATING TOKENS: ``over the moon", ``than meets the eye".
SURROUNDING TOKENS: No interesting patterns.

Step 1.

- The activating tokens are all parts of common idioms.
- The surrounding tokens have nothing in common.

Step 2.

- The examples contain common idioms.
- In some examples, the activating tokens are followed by an exclamation mark.

Step 3.

756 - The activation values are the highest for the more common
 757 idioms in examples 1 and 3.
 758
 759 Let me think carefully. Did I miss any patterns in the text
 760 examples? Are there any more linguistic similarities?
 761
 762 - Yes, I missed one: The text examples all convey positive
 763 sentiment.

764 Afterwards an interpretation is added to the example

765
 766 [interpretation]: Common idioms in text conveying positive
 767 sentiment.

768 A.3 EXAMPLES OF ACTIVATING CONTEXTS AND INTERPRETATIONS

769 As discussed in the main text, the interpretations found for a given latent can be very different
 770 depending on the way to sample the activating contexts shown to the explainer model, see Figure
 771 A2. This has its advantages and disadvantages.
 772
 773

774 When using the top activating contexts, the explainer model normally gives an interpretation that is
 775 more narrow - "The concept of a buffer, referring to something that separates, shields, or protects
 776 one thing from another, often used in various contexts such as physical barriers, chemical reactions,
 777 or digital data processing" - instead of one that captures the full distribution - "Words or phrases
 778 associated with concepts of spatial or temporal separation (buffers, zones, or cushions) or colloqui-
 779 alisms (buff, buffs, or Buffy, referring to a popular TV show or enthusiast)". Here the narrower
 780 interpretations resulted in lower scores: the interpretation generated from examples sampled from
 781 the whole distribution scored 0.95 accuracy in fuzzing and 0.93 accuracy in detection, while the
 782 interpretation generated from top examples only achieves 0.8 accuracy in fuzzing and 0.74 accuracy
 783 in detection.

784 On the other hand, if we look at the second example, sampling from the whole distribution exam-
 785 ples may sometimes confuse the explainer model - "Tokens often precede or succeed prepositions,
 786 articles, and words that signal possession or quantity, frequently indicating a relationship between
 787 objects or actions. of latent", which might not see the pattern that is more clear in the top activating
 788 examples - "Descriptions of food or events where food is involved, often mentioning leftovers, and
 789 sometimes mentioning the act of eating, serving, or storing food, as well as the amount of food,
 790 or the pleasure or satisfaction derived from it." Here we see very poor scores on the interpretation
 791 from randomly sampled examples, 0.54 accuracy in detection and 0.57 accuracy in fuzzing, while
 792 the interpretation generated from the top examples has 0.89 accuracy both in detection and fuzzing,

793 A.4 DIFFERENT SCORING METHODS

794 A.4.1 DETECTION DETAILS

795 The prompt used for detection scoring is the following:

796
 797 You are an intelligent and meticulous linguistics researcher.
 798
 799 You will be given a certain latent of text, such as
 800 ``male pronouns" or ``text with negative sentiment".
 801
 802 You will then be given several text examples. Your task
 803 is to determine which examples possess the latent.
 804
 805 For each example in turn, return 1 if the sentence is
 806 correctly labeled or 0 if the tokens are mislabeled. You must
 807 return your response in a valid Python list. Do not return
 808 anything else besides a Python list.
 809

Together with the prompt, there are several few shot examples like the following:

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

Random

Token Activations: 0.0 0.2 0.5 0.7 1.0

It just shows that he's an ignorant buff
 buffers used with the DRAMs will depend on the
 on the site as a buffer against English-dominated
 is not necessarily coincident with transitional screening
 two buffers.\nSECTION A.14 Vide
 ATED BY DIFFUSIBLE FACTORS\nBUFFINGER
 buffets —.\nTis the season
 skate, kayaking, etc.\nMany buffet style
 used for the buffer nw manager registers shown
 Social Media: Wave of the Now Buffer Been
 league in goal because Gianluigi Buffon
 to womanhood would be buffeted by catacly
 amendment requiring buffers vote, these two i
 data transfer through the buffers is controlled by a handshake
 remain in the feeder pattern of that first school
 local restaurants offering special Easter menus and buffets on
 existence; council rebuffs calls from Mayor Bill de
 which is convertible to restricted shares of Rick
 colored cocker spaniel named Pumpkin graced our family.
 buffer zone between Syria and Jordan which would be
 campaign would cushion the blow of the impending retirement
 25-foot buffer strips on each side of
 Buffy snap-kicked her in the chin.
 Attention Civil War Buffs! Join the Civil War
 trying to send data to the buffer will be halted
 Happy Days. Joan of Arcadia. Buffy the
 really buff (SLIDESHOW)\nMaking the
 observing the "Buffett Rule" that millionaires should
 but as a buffoon.\nSeptember 5
 your tampons...\nFeRD: See
 so maybe he and Buffy could still make it
 could not find anything useful except for a few tamp
 buffer, despite mistreatment—despite some rumors to

Top

Token Activations: 0.0 0.2 0.5 0.8 1.0

design, it is referred to as a buffer bike
 and the angle parking in the downtown provides a buffer
 hectares and a buffer zone of almost 70
 for buffer zone kids to be placed once and their
 which are likely related to its ability to buffer
 in Lyme disease. When buffering species like opossum
 I don't have room in the buffer of
 buffer acids.) When muscles fatigue, acids build-
 162 by a buffer layer
 buffer to at the Big East Media
 25-foot buffer strips on each side of
 the frame to allow the bird some buffer room.
 on the site as a buffer against English-dominated
 nz offers a bigger buffer here (days perhaps)
 and try to build up buffer states between Russia and
 Tests – packaging and buffer vials....
 amendment that a buffer yard be installed on all sides
 space/and/or landscape plantings. A buffer
 does, sometimes buffering agents such as citric acid are
 between uses. A landscaped buffer may be an area
 a powerful H+ buffer inside our muscles called ca
 streaming except for MLB.tv. It buffers all
 Christianity by buffering Europe from Islam, which allows you
 blood volume and, thus, the buffering capacity of
 buffer zone between Syria and Jordan which would be
 amendment requiring buffers vote, these two
 So, the buffered: Planned state in Kanban is
 the blood can buffer, performance declines rapidly. Thus
 Top 10 Secret Buffer Features: Supercharge
 improve water\nBUFFERS: Graduated mix of land
 in buffer About the Author: Jacques van Hee
 buffer, despite mistreatment—despite some rumors to
 point the buffer is only a contingency. [link

Random

Token Activations: 0.0 0.2 0.3 0.5 0.6

rather than the end of the day, as some
 sometimes I carry a cup with me to be c
 result, we ended up having a pretty late lunch
 soup, cabbage stir-fry, salad, toast
 28/2 Use Spikes on Thanksgiving
 these eggs when she cooks congee. But
 HATE TDY! Thank you! I once lost
 plain chicken\nS Yogurt and granola\nD crock
 pics. I'm talking t you too
 have available in your pantry. When we've
 beans or lentils. Served it with brown rice and
 not sure if scalloped potatoes would have made
 not so fragile after refrigeration. In fact since
 short, my mom had about 2 pounds i
 It's all about planning and having the
 No, it's not beef brisket-!
 I'll make more. You're not going
 of Manha for one's family? Obviously
 was also black & burned. The other food was
 are at work. Greta took most of me
 as well (\$7.50). They were
 I can have it on a salad while everyone stuffs
 Oh man, you wouldn't want to be
 of time, so couples who want to practice this
 adoption. Barra said Allen had left the door
 was because it was the last, though I could
 possible type of work was involved in gathering a portion
 dinner?\nSAFEST: Where would you like
 CIOUS. I made this for guests this weekend
 that you put the nutritional facts)\nReplyDeleteAnonymous
 — or until they run out of food. The
 bunny the other day! Bemo commented:
 bases. Also, I always try to leave

Top

Token Activations: 0.0 0.2 0.5 0.8 1.0

will be too much for the two of us
 there and to take home. Variety of food pleased
 we've still got some left over
 ask for take-home boxes for leftovers
 forgot to mention it does not keep well for the
 kept for a late breakfast the next day. The
 double if you have a large crock. I fill
 if you want to serve this the next day
 and I'm looking forward to another warmed-
 excellent as taste and portions large enough for the family
 or like you I might heat up a couple and
 next day - the the government never likes to
 I can see why you got a repeat request for
 and allow diners to leave with them. However
 a huge hit! He had 3 slices lo
 to take leftovers home but we can
 It's all well and good if they want
 pot pork with sweet potatoes "this is great to
 ve had a doggy bag!) In general, we
 leftovers, its up to the dwarves to rescue
 I had tons left over. I sealed the end
 and immediately wanted to make puffy tacos, and they
 let it cool then skim off the fat before rehe
 is \$8\nA large portion of the
 put it away""It's sugar sweet and
 leftovers went back into the fridge. My cousin
 enough food for seconds and thirds Still, not
 we WILL be having it again! My daughter did
 10 had four servings of this with min
 always invite people over for dinner because I know there
 that they may return to enjoy it at a later
 leftover. He took it and his coffee upstairs with
 had never tried anplant before. John Holt Coone

Figure A2: Activating contexts of latent 209 (top) and 293 (bottom), from layers 8 and 32, respectively, of the 131k latent SAE trained on the residual stream of Gemma 2 9b. The shown examples are similar to the ones given to the explainer model to come up with interpretations.

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

```

<user_prompt>

latent interpretation: Words related to American football
positions, specifically the tight end position.

Text examples:

Example 0:Getty Images Patriots tight end Rob Gronkowski
had his boss

Example 1: names of months used in The Lord of the Rings:
the

Example 2: Media Day 2015 LSU defensive end Isaiah Washington
(94) speaks to the

Example 3: shown, is generally not eligible for ads. For
example, videos about recent tragedies,

Example 4: line, with the left side namely tackle Byron
Bell at tackle and guard Amini

<assistant_response>

[1,0,0,0,1]

```

In this example the `<user_prompt>` and the `<assistant_response>` tags are substituted with the correct instruct format used by the scorer model. In detection scoring, 5 shuffled examples are shown to the model at the same time.

A.4.2 FUZZING DETAILS

The prompt used for fuzzing scoring is the following:

```

You are an intelligent and meticulous linguistics researcher.

You will be given a certain latent of text, such as ``male
pronouns" or ``text with negative sentiment". You will
be given a few examples of text that contain this latent.
Portions of the sentence which strongly represent this
latent are between tokens << and >>.

Some examples might be mislabeled. Your task is to determine
if every single token within << and >> is correctly
labeled. Consider that all provided examples could be correct,
none of the examples could be correct, or a mix. An example
is only correct if every marked token is representative
of the latent

For each example in turn, return 1 if the sentence is correctly
labeled or 0 if the tokens are mislabeled. You must return
your response in a valid Python list. Do not return anything
else besides a Python list.

```

Followed by few-shot examples of the kind:

```

<user_prompt>

```

918 latent interpretation: Words related to American football
 919 positions, specifically the tight end position.
 920

921 Text examples:

922

923 Example 0: Getty Images Patriots<< tight end>> Rob Gronkowski
 924 had his boss

925

926 Example 1: posted You should know this<< about>> offensive
 927 line coaches: they are large, demanding<< men>>

928

929 Example 2: Media Day 2015 LSU<< defensive>> end Isaiah
 930 Washington (94) speaks<< to the>>

931

932 Example 3:<< running backs>>, `` he said. .. Defensive
 933 << end>>

934 Carroll Phillips is improving and his injury is

935

936 Example 4:<< line>>, with the left side namely<< tackle>>
 937 Byron Bell at<< tackle>> and<< guard>> Amini

938 <assistant_response>

939 [1,0,0,1,1]

940

941 In this example, the <user_prompt> and the <assistant_response> tags are substituted with the
 942 correct instruct format used by the scorer model. In fuzzing scoring, 5 shuffled examples are shown
 943 to the model at the same time.

944

945 A.4.3 SURPRISAL DETAILS

946 In surprisal scoring, the cross-entropy loss of the scorer model is computed over the tokens of an
 947 example. For each interpretation, this loss is computed with the interpretation and with a default
 948 interpretation - "Various unrelated sentences," where the examples can either be activating context
 949 or non-activating contexts. In our first approach, Llama 3.1 70b base was used. The prompt starts
 950 with few shot examples like:

951 The following is a description of a certain latent of text
 952 and a list of examples that contain the latent.
 953

954 Description:

955

956 References to the Antichrist, the Apocalypse and conspiracy
 957 theories related to those topics.

958

959 Sentences:

960

961 `` by which he distinguishes Antichrist is, that he would
 962 rob God of his honour and take it to himself, he gives
 963 the leading latent which we ought "

964

965 ``3 begins. And the rise of Antichrist. Get ready with "

966

967 `` would be destroyed. The worlds economy would likely
 968 collapse as a result and could usher in a one world government
 969 movement. I wrote a small 6 page "

970 Description:

971

972 Sentences containing digits forming a four-digit year.
 973
 974 Sentences:
 975 `` 20, 2013 at 7:41 pm Martin Smith "
 976
 977 `` of 2012. In other words, Italy's "
 978
 979 ``end 2012 levels). In the first quarter of 2013, we expect
 980 revenue to be up slightly from the fourth quarter "
 981
 982 Description:
 983
 984 Text related to banking and financial institutions.
 985
 986 Sentences:
 987
 988 ``: He is on the Board of Directors with the Lumbee Bank "
 989
 990 `` refurbishing the Bank's branches. BIP reached 400
 991 thousand users in one year The use of BIP has already
 992 doubled The "
 993
 994 `` the Federal Deposit Insurance Corp. "
 995
 996 Description:
 997
 998 Occurrences of the word 'The' at the beginning of sentence.
 999
 1000 Sentences:
 1001
 1002 ``The Smoking Tire hits the canyons with one of the fastest
 Audi's on the road "
 1003
 1004 ``The Chairman of the ABI "
 1005
 1006 ``The administrative center is the town of Koch. "

1007 With the pairs of losses with the interpretation and with the default interpretation, it is possible to
 1008 compute the decrease in loss caused by having access to the interpretation. It is expected that in
 1009 activating contexts this difference will be greater than in non-activating examples, and the score
 1010 of the interpretation is given by the AUC computed using this loss as a proxy for activating and
 1011 non-activating labels.

1012 A.4.4 EMBEDDING DETAILS

1014 For embedding scoring, we use a small 400M parameter model. We chose a small performant
 1015 model on MTEB (Muennighoff et al., 2022) because we found similar scores when using a larger
 1016 7B parameter model, see fig A3, as this size allowed us to do increase the number of examples used
 1017 in scoring.

1018 A set of activating and randomly selected non-activating examples are embedded using the scorer
 1019 model. Then a "query" instruction is embedded:
 1020

1021 Instruct: Retrieve sentences that could be related to
 1022 the interpretation. Query: $\{\text{interpretation}\}$
 1023

1024 The cosine-similarity between the instruction embedding and the examples is computed and is used
 1025 as a proxy for activating and non-activating labels when computing the AUC, which is the score of
 that interpretation.

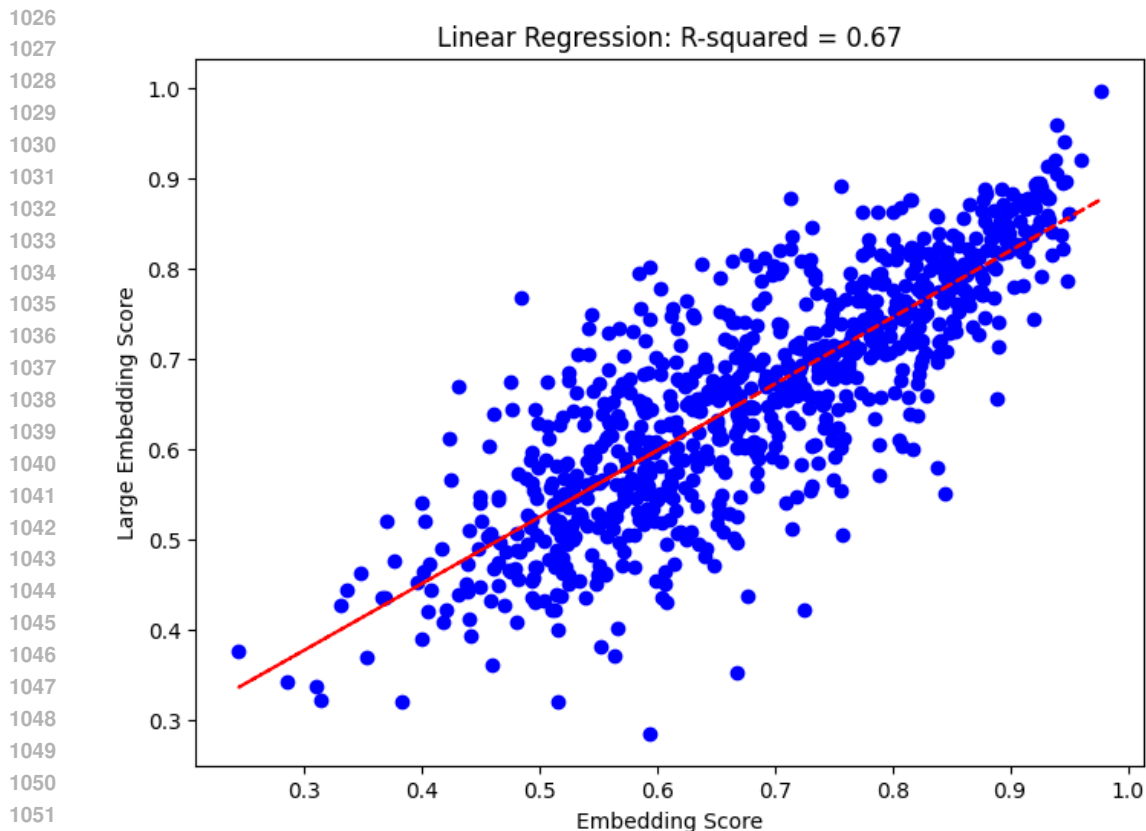


Figure A3: **Comparison between the scores given by a small embedding model and a larger one**

1057 A.4.5 ADVERSARIAL EXAMPLES

1058 Most latents are active in less than 0.1% of the full dataset used to collect the activations, making
1059 random non-activating contexts very diverse. Randomly sampling non-activating examples cannot
1060 be used to determine whether a latent fires in a token in a specific context or on that token in general,
1061 as it is unlikely that that token randomly occurs in each non-activating example. As SAEs are scaled
1062 and latents become sparser and more specific, techniques that overly rely on activating contexts will
1063 have more imprecise results Gao et al. (2024).

1064 Motivated by the phenomenon of latent splitting, we could use “similar” latents to test whether inter-
1065 pretations are precise enough to distinguish between similar contexts. A potential approach is using
1066 cosine similarity between decoder directions of latents to find counterexamples for an interpretation.
1067 Some works (Juang et al., 2024; Macgrath, 2024) have shown that a significant fraction of current
1068 latent interpretations can’t be used to distinguish between latents with high similarity.

1070 A.4.6 CORRELATION BETWEEN SCORES

1071 We compute the correlation between these different scores and simulation 800 different latent inter-
1072 pretations, of the 131k latent SAE trained on the residual stream of Gemma 2 9b spread across 4
1073 different layers. We use 100 activating example, 10 from each of 10 different activating quantiles,
1074 and 100 non-activating examples to compute these scores and use Llama 3.1 70b Instruct as the
1075 scorer model for simulation, fuzzing and detections, while we use Llama 3.1 70b Base for surprisal
1076 and Stella-400M as the embedding model.
1077

1078 We also measure the correlation between our scoring methods and scores given by humans. Humans
1079 were tasked to score how an explanation related to a given context, similar to that done in (Templeton
et al., 2024). The human evaluators could choose between 4 options, ”1 - Explanation not related to

Table A1: Pearson correlation computed over 800 different latent scores

	Fuzzing	Detection	Simulation	Embedding	Surprisal
Fuzzing	1	0.74	0.74	0.42	0.32
Detection		1	0.46	0.70	0.62
Simulation			1	0.30	0.14
Embedding				1	0.79
Surprisal					1

the text”, ”2 - Explanation vaguely related to the text”, ”3 - Incomplete explanation but present in the text”, ”4 - Explanation describes the activation of the tokens”. We measure the average score of a latent in at least 5 contexts. In total 700 contexts and 81 latent interpretation. Because the sample size is smaller, we report the correlation separately from tables A1 and 2, which have significantly better statistics. We find that fuzzing as the best Spearman correlation with human scores - 0.69 - followed by simulation - 0.60 - and recall - 0.59. Surprisal and embedding have 0.34 and 0.32 Spearman correlation respectively.

Interestingly, we also find that the correlation between simulation scores and other scores also increases when the simulation scorer is Claude Sonnet 3.5. These estimations are on a smaller number of latents (100), due to the high cost of simulation using Claude Sonnet, but on this set, the spearman correlation between fuzzing and simulation increases from 0.69 to 0.75, between detection and simulation increases from 0.31 to 0.38, between embedding and simulation from 0.14 to 0.21 and from surprisal to simulation from 0.02 to 0.10. The baseline numbers for the correlation of simulation scores are lower on this set, so we expect that had we done Claude scores on a larger set of explanations, the correlations between the scores would be higher.

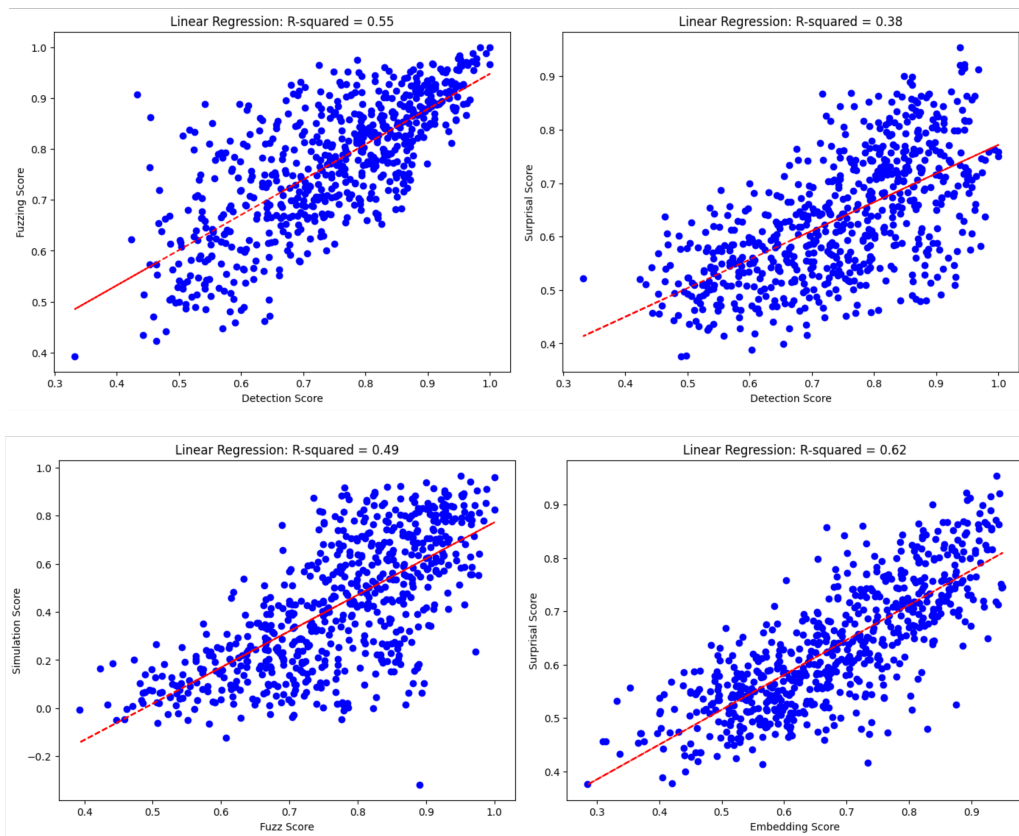


Figure A4: Scatter plots with different combinations of scores

1134 A.4.7 DISAGREEMENT BETWEEN SCORES

1135
1136 In this section, we will discuss how different types of scores might evaluate explanations differently.
1137 In particular, we will focus on an interpretation with fuzzing score but low detection score, high
1138 fuzzing score but low simulation score, high simulation score but low embedding score and high
1139 embedding score and low simulation score.

1140 For instance, feature 2 of layer 8 of the 131k latent SAE model trained on the residual stream of
1141 Gemma 2 9b has a fuzzing score of 0.9 but a detection score of only 0.43. The interpretation given
1142 by our pipeline is "Verbs that link a subject to additional information, often in a formal or descrip-
1143 tive tone." Most of the activations of this latent are on sentences like "This <is needed> because
1144 Magento (...)" or "For instance, this< is> technically correct syntax", where <> represent active
1145 tokens. In both these examples, the model correctly identifies the highlighted tokens as active dur-
1146 ing fuzzing, but incorrectly identifies the context as non-active during detection. On the other hand,
1147 detection incorrectly identifies sentences like "pay for things that would prevent larger issues down
1148 the road is better in the long run." or "Neuroscientist Jack Gallant calls the research a technologic
1149 tour de force and says the ultimate decoder would provide vivid" as active. This explanation is too
1150 vague on the activating context, leading to a low detection score, but specific enough in the types of
1151 tokens that are active, leading to a high fuzz score.

1152 Feature 281 of the same layer has a fuzzing score of 0.97, but a simulation score of only 0.19. Its
1153 interpretation is "URLs or hyperlinks containing query parameters, indicating a request for specific
1154 data or actions on a web page." The simulation score is lower for than the fuzzing score because the
1155 model has to decide which parts of the URL the latent is active on, while the fuzzing scorer either
1156 is shown links that are highlighted which it will say are active or non-links that are highlighted and
1157 are, most likely, not active.

1158 Feature 8 of layer 24 has a simulation score of 0.73, but an embedding score of only 0.43. Its
1159 interpretation is "Prevalence of logical operators and conjunctions in text, including simple addition
1160 and conjunction in various contexts, as well as indicators of contrasting or additional information,
1161 often used for comparison or to provide supplementary details". While doing simulation, the model
1162 correctly identifies that in the sentence "the cost of the unit itself, plus installation and construction
1163 costs.", the token "plus" is active but all the others aren't. The same is true for sentences like "Plus
1164 my potato bread I made on the ANZAC Day" and "Cliff Chiang on reintroducing Orion in Wonder
1165 Woman! Plus interviews". On the other hand, these sentences don't have high embedding similarity
1166 with the explanation.

1167 Feature 10 of layer 16 has the opposite situation happening, with an embedding score of 0.74 and
1168 a simulation score of only 0.04. Its interpretation is "Verbs, prepositions, and adverbs that connect
1169 clauses or indicate direction, movement, or progression, often in a sequence of actions or events".
1170 In sentences like "right, three months ;to close down;. There was dead silence when the message
1171 was read. Everybody waited for Mr. Smith to speak. Mr. Gingham" the simulator model incorrectly
1172 identifies "waited" and the second "to" as active, and misses the real activation in "to close down".
1173 The same happens in "This is ;going; to ;be; one ;swinging; party. Especially if the guests survive
1174 Uncle Wilhelm. I love weddings. I love going to them, I love being in", where the simulator model
1175 identified "love being in" as well as "love going to" as being active. The embedding model correctly
1176 identified that most activating sentences have the mention of motion.

1177 With these examples, we hope to demonstrate that these generated explanations are not perfect,
1178 but that we have an easier understanding of their flaws by looking at which cases different scores
1179 disagree on.

1180 A.5 SAE MODELS USED

1181
1182 Throughout this work, we used different SAEs trained on Gemma. The 16k latent ones trained on
1183 the MLP have the following average L0 norms per layer:

1184	0:50, 1:56, 2:33, 3:55, 4:66, 5:46, 6:46, 7:47, 8:55, 9:40, 10:49, 11:34, 12:42, 13:40, 14:41,
1185	15:45, 16:37, 17:41, 18:36, 19:38, 20:41, 21:34, 22:34, 23:73, 24:32, 25:72, 26:57, 27:52,
1186	28:50, 29:49, 30:51, 31:43, 32:44, 33:48, 34:47, 35:46, 36:47, 37:53, 38:45, 39:43, 40:37,
1187	41:58

Table A2: Impact of training dataset on Fuzzing and Detection performance. Numbers shown are the median score and the interquartile (25%-75%) range.

Experiment	Fuzzing	Detection	Embedding
Random interpretation	0.51 (0.45–0.57)	0.51 (0.45–0.58)	0.51 (0.44–0.57)
Randomly initialized Topk SAE	0.55 (0.50–0.60)	0.54 (0.50–0.59)	–
RPJ-v2	0.76 (0.67–0.86)	0.74 (0.63–0.85)	0.67 (0.57–0.80)
Pile	0.76 (0.67–0.86)	0.76 (0.67–0.85)	0.69 (0.57–0.80)

Table A3: Impact of prompt content on Fuzzing and Detection performance. Numbers shown are the median score and the interquartile (25%-75%) range.

Experiment	Fuzzing	Detection	Embedding
Random interpretation	0.51 (0.45–0.57)	0.51 (0.45–0.58)	0.51 (0.44–0.57)
Randomly initialized Topk SAE	0.55 (0.50–0.60)	0.54 (0.50–0.59)	–
Activations in prompt	0.76 (0.67–0.86)	0.74 (0.63–0.85)	0.68 (0.57–0.80)
No activations in prompt	0.75 (0.65–0.86)	0.73 (0.60–0.84)	0.68 (0.58–0.79)
COT in prompt	0.76 (0.67–0.86)	0.73 (0.61–0.85)	0.65 (0.55–0.75)

The 16k latent ones trained on the residual stream have the following average L0 norms:

0:35, 1:69, 2:67, 3:37, 4:37, 5:37, 6:47, 7:46, 8:51, 9:51, 10:57, 11:32, 12:33, 13:34, 14:35, 15:34, 16:39, 17:38, 18:37, 19:35, 20: 36, 21:36, 22: 35, 23: 35, 24: 34, 25: 34, 26: 35, 27:36, 28: 37, 29:38, 30:37, 31:35, 32: 34, 33:34, 34:34, 35:34, 36:34, 37:34, 38:34, 39:34, 40:32, 41:52

The 131k latent ones trained in the residual steam have the following average L0 norms:

0:30, 1:33, 2:36, 3:46, 4:51, 5:51, 6:66, 7:38, 8:41, 9:42, 10:47, 11:49, 12:52, 13:30, 14:56, 15:55, 16:35, 17:35, 18:34, 19:32, 20:34, 21:33, 22:32, 23:32, 24: 55, 25:54, 26:32, 27:33, 28: 32, 29:33, 30:32, 31:52, 32: 51, 33:51, 34:51, 35:51, 36: 51, 37:53, 38:53, 39:54, 40: 49, 41:45

We also used SAEs with 65k latents, trained on the residual stream and the MLP of Llama 8b.

A.6 FACTORS THAT INFLUENCE THE EXPLAINABILITY

A.6.1 DEPENDENCE ON DATASET

Even though a significant portion of SAE latents are less active when using RPJv2 instead of the Pile, we find that the latents that are active are generally interpretable to the same degree. These evaluations were done using the 131k latent SAE trained on the residual stream of Gemma 2 9b. The scorer and the explainer model were Llama 3.1b 70b instruct, quantized to 4bit.

A.6.2 DEPENDENCE ON CHAIN OF THOUGHT AND ACTIVATION INFORMATION

We find that COT slightly increases the scores of the interpretations found, and that it significantly slows down the rate at which one can produce interpretations. Giving the explainer model, the activations associated with each token seem to slightly increase the scores of the interpretations generated. These evaluations were done using the 131k latent SAE trained on the residual stream of Gemma 2 9b. The scorer and the explainer model were Llama 3.1b 70b instruct, quantized to 4bit.

A.6.3 DEPENDENCE ON CONTEXT LENGTH

The context length of the shown examples did not significantly change the scores obtained for the interpretations. These evaluations were done using the 131k latent SAE trained on the residual stream of Gemma 2 9b. The scorer and the explainer model were Llama 3.1b 70b instruct, quantized to 4 bits.

Table A4: Impact of context length on Fuzzing and Detection performance. Numbers shown are the median score and the interquartile (25%-75%) range.

Experiment	Fuzzing	Detection	Embedding
Random interpretation	0.51 (0.45–0.57)	0.51 (0.45–0.58)	0.51 (0.44–0.57)
Randomly initialized Topk SAE	0.55 (0.50–0.60)	0.54 (0.50–0.59)	–
16 context	0.75 (0.65–0.86)	0.74 (0.62–0.85)	0.70 (0.59–0.81)
32 context	0.76 (0.67–0.86)	0.74 (0.63–0.85)	0.67 (0.57–0.80)
64 context	0.74 (0.64–0.64)	0.70 (0.57–0.81)	0.65 (0.54–0.78)

Table A5: Impact of example sampling strategies on Fuzzing and Detection performance. Numbers shown are the median score and the interquartile (25%-75%) range.

Experiment	Fuzzing	Detection	Embedding
Random interpretation	0.51 (0.45–0.57)	0.51 (0.45–0.58)	0.51 (0.44–0.57)
Randomly initialized Topk SAE	0.55 (0.50–0.60)	0.54 (0.50–0.59)	–
Randomly sampled	0.76 (0.68–0.86)	0.74 (0.62–0.84)	0.66 (0.56–0.78)
Sampled from quantiles	0.77 (0.69–0.87)	0.74 (0.64–0.85)	0.68 (0.57–0.80)
Sampled from top examples	0.73 (0.64–0.83)	0.72 (0.62–0.82)	0.70 (0.58–0.80)

A.6.4 DEPENDENCE ON THE ORIGIN OF EXAMPLES

Sampling the examples from the whole distribution of activations, or sampling a fixed number of examples over different activation quantiles, significantly improved the scores, compared with sampling only from the top examples. These evaluations were done using the 131k latent SAE trained on the residual stream of Gemma 2 9b. The scorer and the explainer model were Llama 3.1b 70b instruct, quantized to 4 bits.

A.6.5 DEPENDENCE ON THE NUMBER OF EXAMPLES

The number of examples shown to the model seems to saturate, at least with the explainer model we used. These evaluations were done using the 131k latent SAE trained on the residual stream of Gemma 2 9b. The scorer and the explainer model were Llama 3.1b 70b instruct, quantized to 4bit.

A.6.6 EXPLAINABILITY ACROSS LAYERS

We find that, in the case of the 131k latent SAE trained on the residual stream, the earliest layers have lower scores than the later layers. These evaluations were done using the 131k latent SAE trained on the residual stream of Gemma 2 9b. The scorer and the explainer model were Llama 3.1b 70b instruct, quantized to 4bit.

Table A6: Impact of number of examples on Fuzzing and Detection performance. Numbers shown are the median score and the interquartile (25%-75%) range.

Experiment	Fuzzing	Detection	Embedding
Random interpretation	0.51 (0.45–0.57)	0.51 (0.45–0.58)	0.51 (0.44–0.57)
Randomly initialized Topk SAE	0.55 (0.50–0.60)	0.54 (0.50–0.59)	–
Shown 10 examples	0.73 (0.62–0.85)	0.71 (0.58–0.82)	0.64 (0.54–0.74)
Shown 20 examples	0.74 (0.64–0.85)	0.72 (0.60–0.84)	0.66 (0.54–0.76)
Shown 40 examples	0.76 (0.67–0.86)	0.74 (0.63–0.85)	0.68 (0.57–0.80)
Shown 60 examples	0.75 (0.66–0.85)	0.73 (0.62–0.84)	0.68 (0.57–0.79)

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

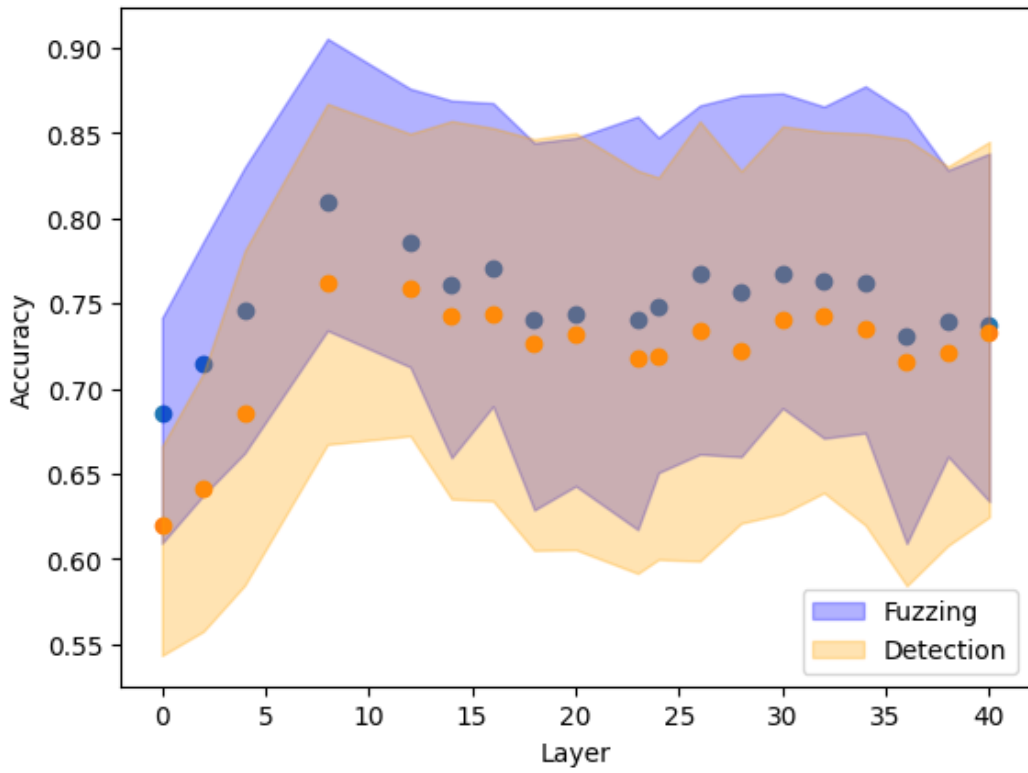


Figure A5: **Accuracy on fuzzing and detection scoring** Dots correspond to the median score over c.a. 300 latent interpretations, and colored region denotes the interquartile range.

Table A7: Comparison of scores of interpretations given by different models. Numbers shown are the median score and the interquartile (25%-75%) range.

Experiment	Fuzzing	Detection	Embedding	Simulation
Random interpretation	0.51 (0.45–0.57)	0.51 (0.45–0.58)	0.51 (0.44–0.57)	-0.02 (-0.02–0.00)
Claude	0.75 (0.68–0.84)	0.75 (0.65–0.85)	0.70 (0.58–0.81)	0.30 (0.28–0.32)
Llama 70b	0.76 (0.67–0.86)	0.74 (0.63–0.85)	0.67 (0.57–0.80)	0.29 (0.28–0.32)
Llama 8b	0.70 (0.60–0.81)	0.70 (0.59–0.81)	0.64 (0.54–0.75)	0.26 (0.23–0.30)
Human	0.75 (0.66–0.85)	0.74 (0.64–0.85)	0.71 (0.62–0.81)	0.36 (0.32–0.39)

Table A8: Comparison of scorer models for fuzzing, detection and simulation. Numbers shown are the median score and the interquartile (25%-75%) range.

Experiment	Fuzzing	Detection	Simulation
Claude	0.76 (0.66–0.88)	0.72 (0.58–0.87)	0.33 (0.29–0.36)
Llama 70b	0.76 (0.67–0.86)	0.74 (0.63–0.85)	0.29 (0.28–0.32)
Llama 8b	0.70 (0.60–0.81)	0.70 (0.59–0.81)	0.26 (0.25–0.31)

A.6.7 DEPENDENCE ON SIZE OF EXPLAINER MODELS

Having a larger explainer model leads to better interpretations, but our results suggest that this benefit tends to saturate, as we find that the scores of explanations given by Claude or Llama 70b are very similar. We only did human interpretation of 100 latents. Simulations scores are also only computed over 100 latents.

A.6.8 DEPENDENCE ON SIZE OF EXPLAINER MODELS

We find that both fuzzing, detection and simulation scoring are dependent on the size of the explainer model, that is, given the same explanations, larger models achieve higher scores on average. While Claude and Llama 80b have similar fuzzing and detection scores, simulation, being more complex, is better performed by Claude.

A.6.9 DEPENDENCE ON SAE SIZE AND LOCATION

We compare the interpretability of SAEs with different number of latents, trained on the same model, with the neurons of that model. We find that the SAE with the highest number of latents to have the highest scores in the case of the Gemma 2 9b SAEs and that SAEs trained on the residual stream have higher scores for both Gemma 2 9b and Llama 3.1 8b.

The scorer and the explainer model were Llama 3.1 70b instruct, quantized to 4bit.

A.7 AUTOMATICALLY INTERPRETING INTERVENTIONS

A.7.1 SCORING IMPLEMENTATION

For each latent, we sample a pool of 40 length-64 prompts from RPJv2 (Computer, 2023), of which 30 i.i.d. prompts are taken for scoring, while the remaining 10 are used by the explainer. The pool is sampled among nonzero activating contexts, stratified by the quintile of the latent’s max activation in that context (i.e., the nonzero activating documents are sorted by the context’s maximum activation of the current latent, then 8 contexts are sampled from the first quintile, 8 from the second etc.). Each of these prompts is then truncated to include only the first token on which the latent activates, so that all previous activations for the latent are 0. Activations on the first token position are ignored because the SAEs were not trained on that position.

We filter out latents that activates on less than 200 of the 10 million RPJ-v2 tokens.

We then sample generations with a maximum of 8 new tokens from the subject model. For generations with the intervention, we perform an additive intervention on the latent at all token positions

Table A9: Comparison of SAEs for Fuzzing and Detection. Numbers shown are the median score and the interquartile (25%-75%) range.

Experiment	Fuzzing	Detection
Random interpretation	0.51 (0.45–0.57)	0.51 (0.45–0.58)
Randomly initialized Topk SAE	0.55 (0.50–0.60)	0.54 (0.50–0.59)
131k latents Gemma 2 9b	0.76 (0.67–0.86)	0.74 (0.63–0.85)
16k latents Gemma 2 9b	0.73 (0.63–0.83)	0.70 (0.59–0.79)
Top 32 neurons Gemma 2 9b	0.62 (0.54–0.70)	0.59 (0.53–0.64)
Top 256 neurons Gemma 2 9b	0.59 (0.53–0.65)	0.57 (0.52–0.62)
262k latents Llama 3.1 8b MLP	0.79 (0.69–0.86)	0.79 (0.64–0.85)
262k latents Llama 3.1 8b	0.81 (0.71–0.86)	0.83 (0.68–0.85)
Top 32 neurons Llama 3.1 8b	0.55 (0.51–0.62)	0.53 (0.50–0.59)
Top 256 neurons Llama 3.1 8b	0.54 (0.49–0.60)	0.53 (0.50–0.57)

after and including the final prompt token. We tune the intervention strength of each latent to various KL-divergence values on the scoring set with a binary search. We stop the binary search when the KL divergence is within 10% of the desired value³. For our zero-ablation experiment, instead of doing an additive intervention we clamp the latent’s activation to 0. Specifically, the hidden states are encoded by the SAE, then the SAE reconstruction error is computed using a clean decoding, then the SAE encoding is clamped and decoded, and the error is added to the clamped decoding.

The scorer is Llama-3.1 8B base. We use the base model for improved calibration, and prompt it as follows.

```

<PASSAGE>
from west to east, the westmost of the seven wonders of the
world is the great wall of china

The above passage contains an amplified amount of "Asia"

<PASSAGE>
Given 4x is less than 10, 4

The above passage contains an amplified amount of "numbers"

<PASSAGE>
In information theory, the information content, self-
information, surprisal, or Shannon information is a basic
quantity derived by her when she was a student at Windsor

The above passage contains an amplified amount of
"she/her pronouns"

<PASSAGE>
My favorite food is oranges

The above passage contains an amplified amount of
"fruits and vegetables"

<PASSAGE>
...
```

³Sometimes the KL-divergence is not perfectly monotonic in the intervention strength so 10% error is exceeded. We report the average KL divergence that we observe in Figure A6

1458 A.7.2 GENERATING INTERPRETATIONS

1459 We generate interpretations using only the intervention’s effect on the subject’s next-token probabilities because this leads to a concise and precise prompt. The explainer, like the scorer, is Llama-3.1
1461 8B.

1463 The explainer sees a distribution of 10 prompts that is sampled i.i.d. from the same population as
1464 the scorer’s prompts.

1465 We use the following prompt for the explainer, with 3 few-shot examples truncated for brevity.
1466

1467 We’re studying neurons in a transformer model. We want to know
1468 how intervening on them affects the model’s output.

1469 For each neuron, we’ll show you a few prompts where we
1470 intervened on that neuron at the final token position, and the
1471 tokens whose logits increased the most.
1472

1473 The tokens are shown in descending order of their probability
1474 increase, given in parentheses. Your job is to give a short
1475 summary of what outputs the neuron promotes.
1476

1477 Neuron 1
1478 <PROMPT>Given 4x is less than 10,</PROMPT>
1479 Most increased tokens: ‘ 4’ (+0.11), ‘ 10’ (+0.04),
1480 ‘ 40’ (+0.02), ‘ 2’ (+0.01)

1481 <PROMPT>For some reason</PROMPT>
1482 Most increased tokens: ‘ one’ (+0.14), ‘ 1’ (+0.01),
1483 ‘ fr’ (+0.01)
1484

1485 <PROMPT>insurance does not cover claims for accounts
1486 with</PROMPT>
1487 Most increased tokens: ‘ one’ (+0.1), ‘ more’ (+0.02),
1488 ‘ 10’ (+0.01)
1489

1490 interpretation: numbers

1491 Neuron 2

1492 ...
1493

1494 A.7.3 BASELINES

1496 **Random SAE.** We experiment with a random top-k SAE with $k = 50$, roughly the average sparsity
1497 of the gamma SAEs. The encoder is initialized with 131,072 spherically uniform unit-norm latents,
1498 and the decoder is initialized to its transpose. We use a random SAE with TopK activations because
1499 we need sparsity for our sampling procedure to work properly.

1500 **Random interpretations.** For each layer and target KL value, we compute a random interpretation
1501 baseline where we shuffle the interpretations across latents.
1502

1503 A.7.4 INTERVENTION INTERPRETABILITY RESULTS

1505 Figure A6 shows the average intervention score at each layer, for a few KL values. Intervention
1506 scores are higher in later layers, likely because of the increased proximity to the model’s output.

1507 Figure A7 is similar to the right half of Figure 4, but at multiple layers and KL values. Early layers
1508 have small intervention scores across the board, so there is little correlation, while the reason for
1509 little correlation in layer 41 is less clear.
1510
1511

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

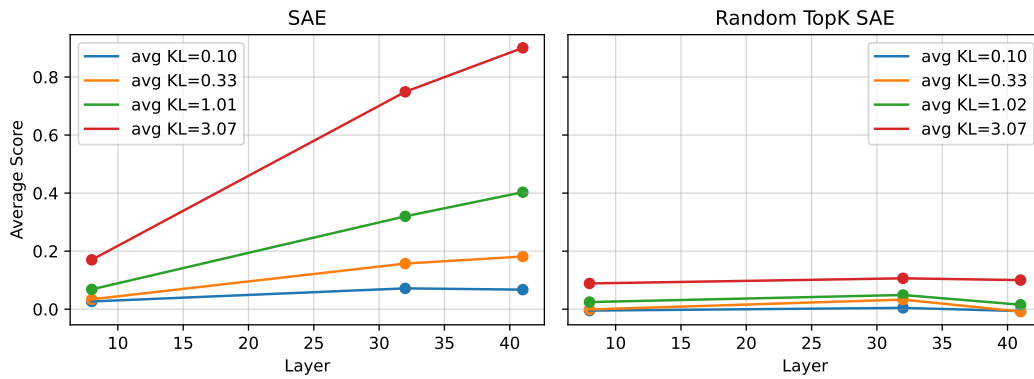


Figure A6: Average intervention score vs layer. SAE latents in later layers have more explainable effects on output. Random latents, however, have uninterpretable effects on output, even at late layers.

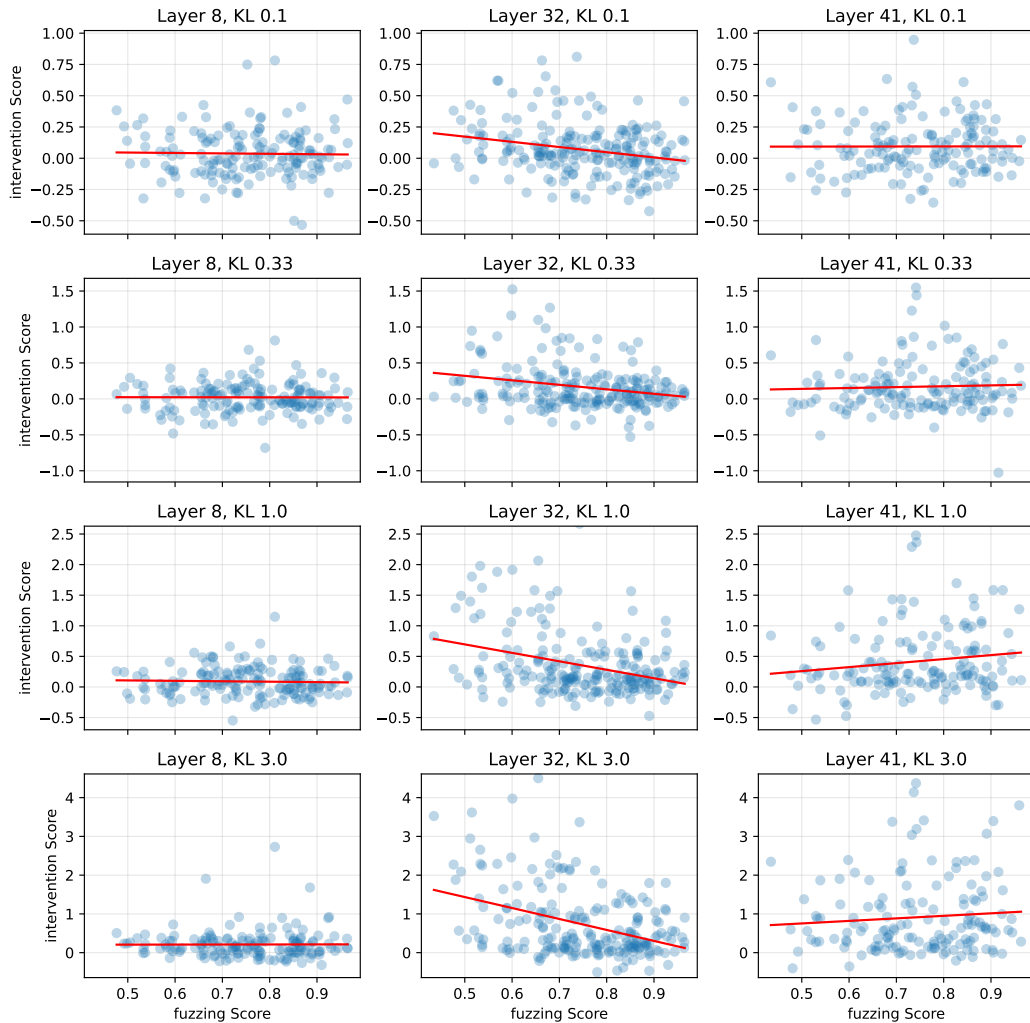


Figure A7: Comparison of intervention and fuzzing scores across layers at various target KL values.