

AUTOMATED INTERPRETABILITY METRICS DO NOT DISTINGUISH TRAINED AND RANDOM TRANSFORMERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Sparse autoencoders (SAEs) are widely used to extract sparse, interpretable latents from transformer activations. We test whether commonly used SAE quality metrics and automatic explanation pipelines can distinguish trained transformers from randomly initialized ones (e.g., where parameters are sampled i.i.d. from a Gaussian). Over a wide range of Pythia model sizes and multiple randomization schemes, we find that, in many settings, SAEs trained on randomly initialized transformers produce auto-interpretability scores and reconstruction metrics that are similar to those from trained models. These results show that high aggregate auto-interpretability scores do not, by themselves, guarantee that learned, computationally relevant features have been recovered. We therefore recommend treating common SAE metrics as useful but insufficient proxies for mechanistic interpretability and argue for routine randomized baselines and targeted measures of feature ‘abstractness.’

1 INTRODUCTION

Sparse autoencoders (SAEs) are a popular tool in mechanistic interpretability research, with the aim of disentangling the internal representations of neural networks by learning sparse, interpretable features from network activations (Elhage et al., 2022; Sharkey et al., 2022; Cunningham et al., 2023; Bricken et al., 2023). An autoencoder with a high-dimensional hidden layer is trained to reconstruct activations while enforcing sparsity (Gao et al., 2024; Templeton et al., 2024; Lieberum et al., 2024), with the aim of discovering the underlying concepts or ‘features’ learned by the network (Park et al., 2023; Wattenberg and Viégas, 2024). Developing better SAEs relies on quantitative evaluation metrics like auto-interpretability scores that measure agreement between generated explanations and activation patterns (Bills et al., 2023; Paulo et al., 2024; Karvonen et al., 2024a).

For an interpretability method to be considered robust, its evaluation metrics should distinguish features learned through training from artifacts arising from the data or model architecture. A key sanity check is therefore to compare the method’s output on a trained model against a strong null model, such as one with randomly initialized weights (Adebayo et al., 2020). We apply this sanity check to SAEs and find that several common quantitative metrics do not always clearly distinguish between the trained and randomized settings. In particular, we found that SAEs trained on transformers with random parameters can yield latents with auto-interpretability scores (Bills et al., 2023; Paulo et al., 2024) that are surprisingly similar to those from a fully trained model.

This result raises important questions about what we can glean from applying these metrics of SAE quality. High auto-interpretability scores alone do not guarantee that an SAE has identified complex, learned computations. Instead, such scores may sometimes reflect simpler statistical properties of the training data (Dooms and Wilhelm, 2024) or architectural inductive biases that are present even without training. Indeed, one could argue that a randomly initialized network still performs a basic form of computation, such as preserving or amplifying the sparse structure of its inputs (Section 4). From this perspective, SAEs might faithfully interpret this simple, inherent computation.

While some SAE features from trained models clearly arise from learned computation, the commonly used aggregate metrics are often insufficient for determining whether a given SAE has learned these more complex features. These results have important implications for mechanistic interpretability research. In particular, we suggest that more rigorous methods to distinguish between artifacts and genuinely learned computations are needed, and that interpretability techniques should be carefully validated against appropriate null models.

Finally, we speculate about why these patterns might emerge. At a high level, there are two hypotheses: (1) the input data already exhibits superposition, and randomly initialized neural networks largely preserve this superposition; and (2) randomly initialized neural networks amplify or even introduce superposed structure to the input data (e.g., given dense input generated i.i.d. from a Gaussian). We present toy models to demonstrate the plausibility of these hypotheses in Section 4 but defer conclusions as to the mechanism responsible to future work.

2 RELATED WORK

Sparse dictionary learning Under a different name, ‘superposition’ in visual data is one of the foundational observations of computational neuroscience. Olshausen and Field (1996; 1997) showed that the receptive fields of simple cells in the mammalian visual cortex can be explained as a result of sparse coding, i.e., representing a relatively large number of signals (sensory information) by simultaneously activating a small number of elements (neurons). Coding theory offers a perspective on efforts to extract the ‘underlying signals’ responsible for neural network activations (Marshall and Kirchner, 2024).

Sparse dictionary learning (SDL) approximates a set of input vectors by linear combinations of a relatively small number of learned basis vectors. The learned basis is usually overcomplete: it has a greater dimension than the inputs. SDL algorithms include Independent Component Analysis (ICA), which finds a linear representation of the data such that the components are maximally statistically independent (Bell and Sejnowski, 1995; Hyvärinen and Oja, 2000). Sparse autoencoders (SAEs) are a simple neural network approach (Lee et al., 2006; Ng, 2011; Makhzani and Frey, 2014). Typically, an autoencoder with a single hidden layer that is many times larger than the input activation vectors is trained with an objective that imposes or incentivizes sparsity in its hidden layer activations to try to find this structure. A **latent** is a single neuron (dimension) in the autoencoder’s hidden layer.

Mechanistic interpretability Recently, it has become common to understand ‘features’ or concepts in language models as low-dimensional subspaces of internal model activations (Park et al., 2023; Wattenberg and Viégas, 2024; Engels et al., 2024). If such sparse or ‘superposed’ structure exists, we expect to be able to ‘intervene on’ or ‘steer’ the activations, i.e., to modify or replace them to express different concepts and so influence model behavior (Meng et al., 2022; Zhang and Nanda, 2023; Heimersheim and Nanda, 2024; Makelov, 2024; O’Brien et al., 2024).

SAEs are a popular approach for discovering features, where one typically trains a single autoencoder to reconstruct the activations of a single neural network layer, e.g., the transformer residual stream (Sharkey et al., 2022; Cunningham et al., 2023; Bricken et al., 2023). Many SAE architectures have been suggested, which commonly vary the activation function applied after the linear encoder (Makhzani and Frey, 2014; Gao et al., 2024; Rajamanoharan et al., 2024b; Lieberum et al., 2024). SAEs have also been trained with different objectives (Braun et al., 2024; Farnik et al., 2025) and applied to multiple layers simultaneously (Yun et al., 2021; Lawson et al., 2024; Lindsey et al., 2024).

Besides reconstruction errors and preservation of the underlying model’s performance, SAEs have been evaluated according to whether they capture specific concepts (Gurnee et al., 2023; Gao et al., 2024) or factual knowledge (Huang et al., 2024; Chaudhary and Geiger, 2024), and whether these can be used to ‘unlearn’ concepts (Karvonen et al., 2024b).

Automatic neuron description SAEs often learn tens of thousands of latents, which are infeasible to describe by hand. Yun et al. (2021) find the tokens that maximally activate a dictionary element from a text dataset and manually inspect activation patterns. Instead, researchers typically collect latent activation patterns over a text dataset and prompt a large language model to explain them (e.g. Bills et al., 2023; Foote et al., 2023). These methods have been widely adopted (e.g. Cunningham et al., 2023; Bricken et al., 2023; Gao et al., 2024; Templeton et al., 2024; Lieberum et al., 2024).

Bills et al. (2023) generate an explanation for the activation patterns of a language-model neuron over examples from a dataset, simulate the patterns based on the explanation, and score the explanation by comparing the observed and simulated activations. This method is commonly known as auto-interpretability (as in self-interpreting). Paulo et al. (2024) introduce classification-based measures of the fidelity of automatic descriptions that are inexpensive to compute relative to simulating activation patterns and an open-source pipeline to compute these measures. Choi et al. (2024) use best-of- k

sampling to generate multiple explanations based on different subsets of the examples that maximally activate a neuron. Importantly, they fine-tune Llama-3.1-8B-Instruct on the top-scoring explanations to obtain inexpensive ‘explainer’ and ‘simulator’ models.

Polysemanticity Lecomte et al. (2024) noted that neurons may become polysemantic incidentally. A polysemantic neuron (basis dimension) of a network layer represents multiple interpretable concepts (Elhage et al., 2022; Scherlis et al., 2023); unsurprisingly, individual neurons in a randomly initialized network may be polysemantic. By contrast, our work studies *superposition* (Elhage et al., 2022; Chan, 2024), which pertains to the representations learned across a whole network layer as opposed to any individual neuron. In particular, superposition allows a network layer as a whole to represent a larger number of (sparse) features than the layer has (dense) neurons by sparse coding (only a few concepts are active at a time, i.e., a given token position).

Training only the embeddings Zhong and Andreas (2024) showed that transformers learn surprising algorithmic capabilities when only the embeddings are trained and no other parameters. These results demonstrate that the behavior of a randomly initialized transformer can be shaped to a surprising extent by training only a few parameters. However, our setting is very different: besides considering SAEs, we randomize *all* the parameters, including the embeddings, in our ‘Step-0’ and ‘Re-randomized incl. embeddings’ variants. Our ‘Re-randomized excl. embeddings’ variant uses pre-trained embeddings, but we do not train those embeddings with fixed, randomized weights. Instead, we freeze the pre-trained embeddings and randomize the other weights (Section 3).

Random transformers for board games Karvonen et al. (2024c) found that SAEs were considerably better at extracting meaningful structure from chess games using pre-trained transformers, as opposed to those with random weights. However, the data from board games is wildly different from language data. In particular, there is reason to expect that language is sparse (e.g., a particular concept such as ‘serendipitous’ appears only rarely), and that this sparse structure is ‘aligned’ with conceptual meaning. In contrast, in board games, this is not necessarily true: a useful concept such as a knight fork does not necessarily turn up sparsely in board games.

Random one-layer transformers Bricken et al. (2023) found that auto-interpretability scores discriminated effectively between random and trained one-layer transformers. Similarly, we found that auto-interpretability scores for randomized models were relatively low for smaller models (e.g., Pythia-70m) but that the gap was narrowed for larger models (e.g., Pythia-6.9b).

3 RESULTS

We trained per-layer SAEs on the residual stream activation vectors of transformer language models from the Pythia suite, with between 70M and 7B parameters (Biderman et al., 2023). We compared SAEs trained on different variants of the underlying transformers:

- **Trained:** The usual, trained model.
- **Re-randomized incl. embeddings:** All the model parameters, including the embeddings, are re-initialized by sampling Gaussian noise with mean and variance equal to the values for each of the original, trained weight matrices.
- **Re-randomized excl. embeddings:** As above, except the embedding and unembedding weight matrices are not re-initialized, i.e., are the same as the original, trained model.
- **Step-0:** For Pythia models, the step0 revisions are available, which are the original model weights at initialization, i.e., before any learning (Biderman et al., 2023).
- **Control:** The original, trained model, except where the input token embeddings are replaced at inference time by sampling i.i.d. standard Gaussian noise for each token, such that a given token does not have a consistent embedding vector. For this variant, we expect auto-interpretability to perform at the level of chance.

For our primary experiments, we trained SAEs on 100M tokens from the RedPajama dataset (Weber et al., 2024) using an activation buffer size of 10M tokens (see Appendix C for a subset of experiments

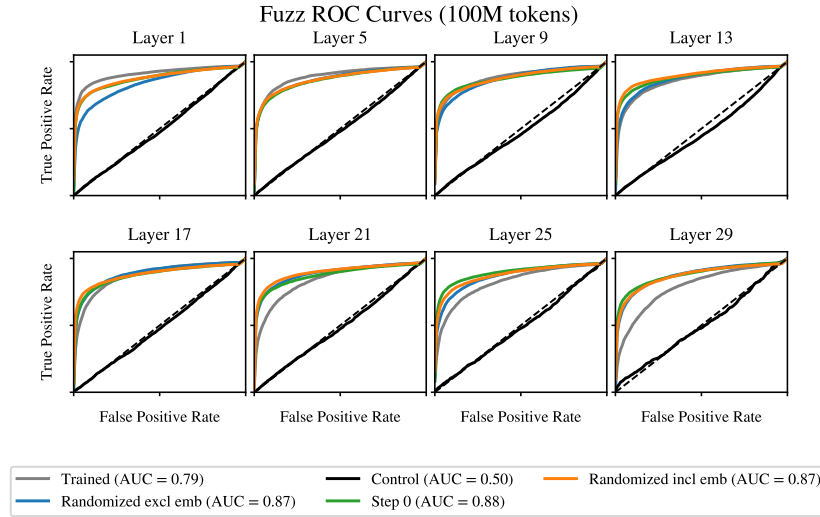


Figure 1: ‘Fuzzing’ ROC curve vs. layer for Pythia-6.9b (100 latents sampled per SAE). The trained model (gray line) and randomized variants (colored) overlap, whereas the control (black) is near chance (dotted). This suggests aggregate AUROC alone is insufficient to attribute latents to learned computation. See Figure 2 for other metrics/model sizes and Appendix E for multiple random seeds.

that demonstrate similar results with SAEs trained on one billion tokens). For models with fewer than 410M parameters, we trained an SAE at every layer; for Pythia-1b, we trained SAEs at every second layer; and for Pythia-6.9b, we trained SAEs at every fourth layer.

Unless otherwise stated, we trained k -sparse autoencoders (also known as TopK SAEs; Makhzani and Frey 2014; Gao et al. 2024), with an expansion factor of $R = 64$ and sparsity $k = 32$. We confirm that our results are robust with respect to these hyperparameters by training SAEs on Pythia-160m with expansion factors equal to powers of 2 between 16 and 128, and sparsities of 16 and 32 (Figure 18). The training implementation is based on Belrose et al. (2025); our evaluations are based on Caden et al. (2025) and Karvonen et al. (2024a).

Auto-interpretability Feature explanations that identify a concept can be input to a classifier that predicts whether the concept appears in the text inputs. Such a classifier may be evaluated by traditional metrics, like the area under the receiver operating characteristic (ROC) curve (AUROC). Paulo et al. (2024) proposed ‘fuzzing’ and ‘detection’ classification tasks to evaluate feature explanations. For ‘fuzzing’ scoring, both positive and negative examples of tokens (i.e., with non-zero and zero activation values, respectively) for a given latent are delimited with special characters, and a language model is prompted to identify which examples have been correctly delimited for the latent given its explanation. For ‘detection’, a language model is asked to identify which examples contain activating tokens for each feature. Bills et al. (2023) originally proposed ‘simulation’ scoring, based on the correlation between predicted and observed activations, but this method is expensive to compute.

Except where noted, we report ‘fuzzing’ scores as a measure of auto-interpretability, because this measure has been demonstrated to correlate with simulation scoring (Paulo et al., 2024). We include similar AUROC curves for the ‘detection’ scoring method in Appendix B. For each trained SAE (i.e., underlying model, variant, and layer), we randomly sampled 100 features to obtain auto-interpretability scores. The implementation is based on Paulo et al. (2024). We use the Meta-Llama-3.1-70B-Instruct-AWQ-INT4 model to generate explanations and make predictions (larger than the 8B models used by Choi et al. (2024) and open-source, unlike Bills et al. 2023).

We found that the auto-interpretability scores were far more similar between the trained and randomized models than with the control (Figures 1 and 2). The similarity between the ROC scores for trained and randomized transformers demonstrates that ‘fuzzing’ auto-interpretability alone, applied to SAE latent explanations, may not meaningfully distinguish between these underlying models.

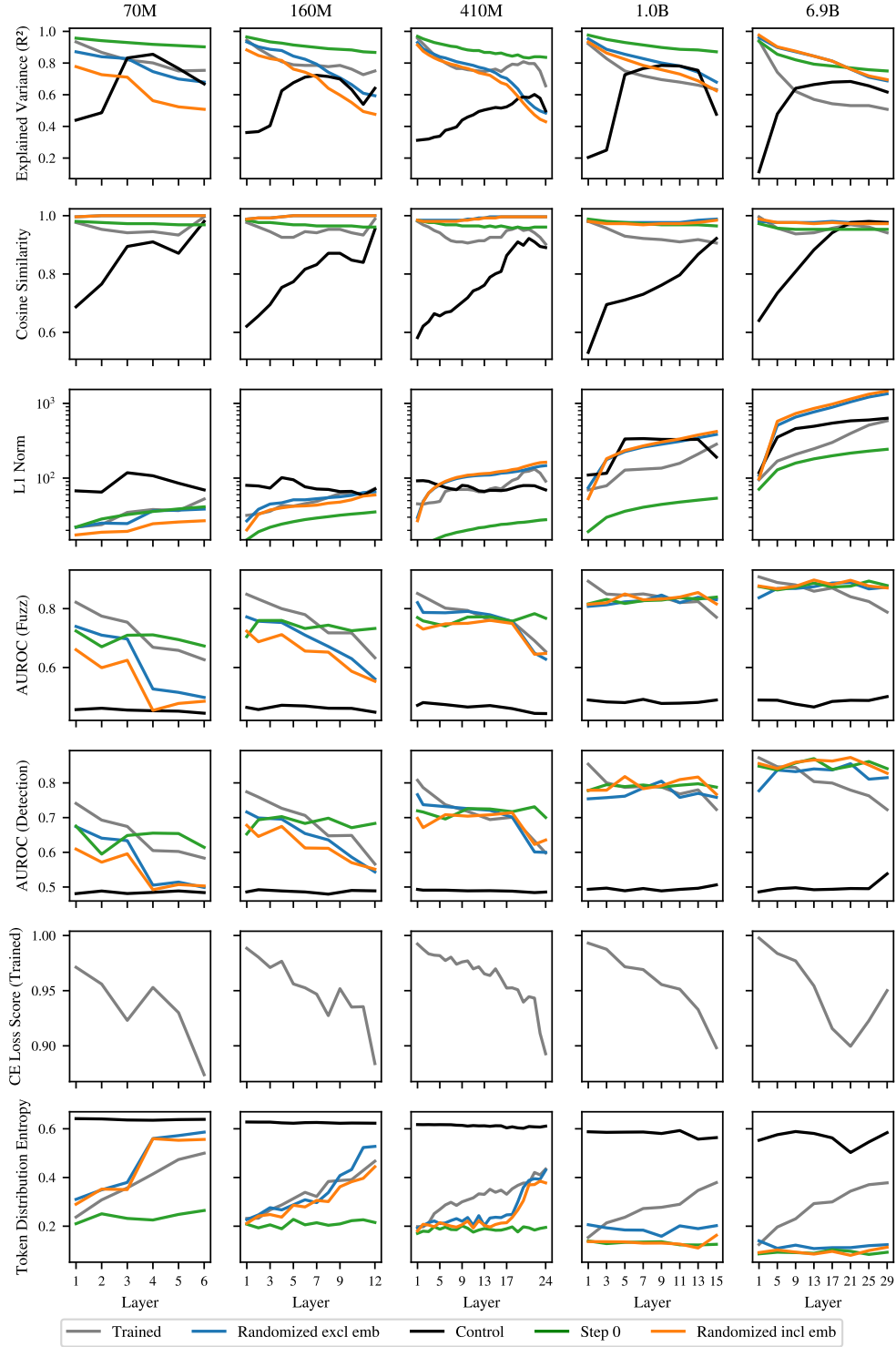


Figure 2: Comparison of sparse autoencoder performance across Pythia models (70M to 6.9B parameters). The different SAE variants show remarkably similar trends across model scales, with larger models exhibiting more consistent behavior across layers. All variants save for control achieve comparable performance despite fundamentally different initialization approaches.

Evaluation We considered standard SAE evaluation metrics alongside the auto-interpretability AUROC for Pythia models with between 70M and 6.9B parameters. As above, we broadly found that the randomly initialized and re-randomized models (Figure 2; blue, green, orange lines) were more similar to the trained model (Figure 2; gray lines) than to our control (Figure 2; black lines).

Notably, the cosine similarity between the original activation and the SAE reconstruction and explained variance are often far lower for the random control than the other models, and its reconstruction errors tend to increase across layers while the remaining variants decrease. For the random control, this can perhaps be explained by the fact that a Gaussian is the highest entropy distribution with fixed mean and variance (Jaynes, 2003); we speculate that Gaussian vectors are the ‘least structured’, in some sense, and thus hardest for SAEs to reconstruct. As Gaussian-distributed activations are propagated through successive layers, we would expect the activations to become less Gaussian and perhaps more ‘sparse’, i.e., easier to reconstruct (Section 4).

Interestingly, the randomized variants (blue and orange lines) are more similar to the trained model than the variant at initialization (green line). This is especially evident if we look at the L^1 norm values in larger models. We speculate that this pattern arises because parameter norms may differ greatly between a trained model and its state at initialization. In contrast, our randomization procedure was specifically designed to preserve parameter norms with respect to the trained model. The scale of parameters at different layers may be important, e.g., to control the growth of activations as they progress through the residual stream (Liu et al., 2020). In the AUROC plots, we find that for all but the control variant, AUROC increases with model size. We speculate that features become more specific as SAE size increases: in smaller SAEs, each latent must explain more of the input, making classification tasks easier for larger SAEs.

Figure 2 (row five) shows the cross-entropy (CE) loss score, or loss recovered, against model layer. This is the increase in the loss when the original model activations are replaced by their SAE reconstructions, divided by the increase when the activations are replaced by zeros (‘ablated’). The results show that the ‘trained’ variant SAEs perform similarly to others from the literature (e.g., Kissane et al., 2024; Rajamanoharan et al., 2024a; Mudide et al., 2024). Importantly, the CE loss score only makes sense for the trained variant: for any of the randomized variants, the loss is very poor, regardless of whether the original or reconstructed activations are used.

Latent explanation complexity Despite sometimes similar auto-interpretability scores and evaluation metrics, we had expected that SAEs applied to trained vs. randomized transformers would discover qualitatively different features. In particular, we expected SAEs trained on the randomized variants to learn relatively simple features based on characteristics of the input text, but not more complex, abstract features as with trained transformers (Templeton et al., 2024). For qualitative examples, we provide a random sample of features and the corresponding maximally activating dataset examples for each variant of Pythia-6.9b in Appendix J, and more detailed information in Appendix L.

Anecdotally, we have observed that a significant proportion of SAE latents have non-zero activations only on a single token or a small number of distinct tokens within a text dataset (e.g., Lin and Bloom, 2024; Dooms and Wilhelm, 2024). Hence, a simple measure of the complexity of an SAE latent given a set of maximally activating examples is the degree to which the latent activates on a single token ID or multiple distinct IDs. Specifically, we quantify the number of token IDs in terms of the entropy of the observed distribution of latent activations over tokens: the greater the entropy, the more ‘spread out’ the latent activations, and the less token-specific the latent. We take this distribution to be the total latent activation per token across the set of maximally activating examples used to generate explanations for auto-interpretability. We show the relationship between entropy and ‘fuzzing’ AUROC score for individual latents in Appendix H.

We include the entropy of the observed distributions of latent activations over token IDs in the last row of Figure 2. The negative control variant displays a consistently high entropy, which is to be expected given that the embedding for a given token ID is sampled i.i.d. from a Gaussian on each occurrence of the token, i.e., a token does not have a consistent embedding vector (Section 3). For the trained variant, the entropy increases across layers, i.e., the further into the model, the less likely the maximally activating examples for each latent contain activations concentrated on a single token. This is also expected: at later layers, we expect more abstract features that are less similar to token

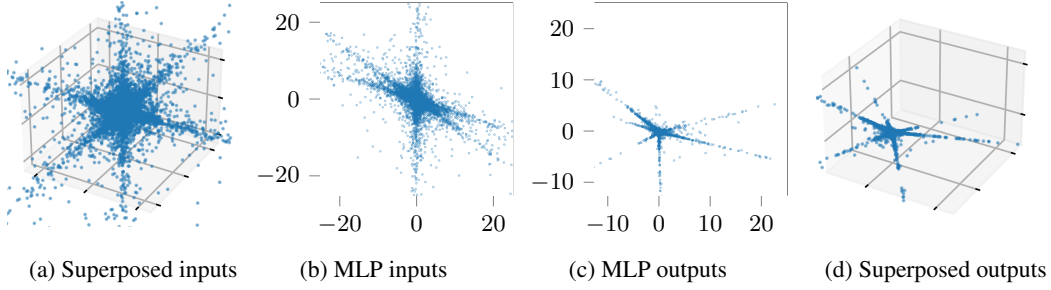


Figure 3: An example of the effect of a randomly initialized neural network on superposed input data. We take 10K samples of $n_s = 3$ sparse input features from a Lomax distribution with shape $\alpha = 1$ and scale $\lambda = 1$ and project these to $n_d = 2$ dense input features by an i.i.d. standard normal matrix. Then, we pass the dense outputs to a two-layer MLP with ReLU activation and hidden size of $4n_d$ and recover $n_s = 3$ sparse outputs by the inverse of the previously generated projection matrix.

embeddings. Finally, the entropy for randomized models tends to be lower than for either the trained or control variants, indicating that latents are activated specifically at one or a few IDs.

In combination with the preceding results, this suggests that standard SAE quality and auto-interpretability metrics are missing an important aspect of SAE features: their ‘abstractness’. While the token distribution entropy is not a direct measure of ‘abstractness’, it suggests that the randomized variants, viewed in the context of their similar auto-interpretability scores to the trained variant, remain able to learn simple, single-token features. However, unlike the trained variant, the features of the randomized variants do not become more complex as the layer index increases.

4 A TOY MODEL OF SUPERPOSITION IN RANDOM NETWORKS

We speculated in Section 1 that the apparently high degree of sparsity and interpretability in the activations of randomized transformers might be because the input data exhibits superposition, which neural networks preserve, or neural networks somehow amplify or even introduce superposition into the input data. In this section, we examine both possibilities through the lens of toy models. We find some evidence to support each potential cause, but we leave the question of which predominates in the case of randomized transformers and the results detailed in the main text to future work.

4.1 MATRIX MULTIPLICATIONS PRESERVE SUPERPOSITION

First, we consider a simplified model to demonstrate that multiplication by a weight matrix W preserves superposition. Imagine that we generate superposed input data x by first generating n_s i.i.d. ‘sparse’ features z from a heavy-tailed Lomax distribution $z \sim \text{Lomax}(\alpha, \lambda)$. We can project the higher-dimensional, sparse z down to lower-dimensional, dense x with a matrix D , then add Gaussian noise with a small variance Σ , $x \sim \mathcal{N}(x; Dz, \Sigma)$. Importantly, if we multiply x by some matrix W , then $x' = Wx$ is *also* superposed: it is generated by the same model as x , except with different noise covariances and mappings from z to x , namely $x' \sim \mathcal{N}(z; W Dz, W \Sigma W^T)$.

We can see that the same intuition might extend to neural networks with nonlinearities by visualizing the results of passing the dense activations through a simple feed-forward network (MLP). Figure 3 shows an example where $n_s = 3$ sparse features are projected down to $n_d = 2$ dense features, and the MLP outputs appear superposed despite the non-linearity. Moreover, it suggests that NNs might amplify superposition rather than only preserving it: comparing the inputs (Figures 3a and 3b) to the outputs (Figures 3d and 3c), there are fewer points between the ‘arms’ of the outputs.

4.2 DO RANDOM NNs PRESERVE OR AMPLIFY SUPERPOSITION?

We investigated this suggestion by generating toy data with the same procedure as Sharkey et al. (2022), i.e., sampling ground-truth features on a hypersphere and generating correlated feature coefficients such that only a small number are active (Appendix I.1). We then passed these inputs

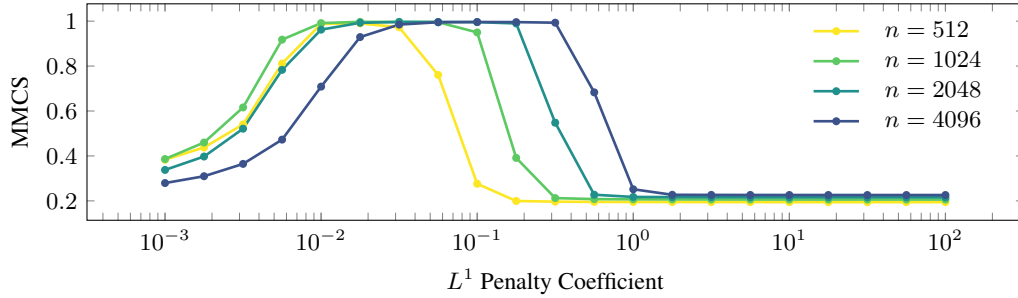


Figure 4: The mean max cosine similarity (MMCS) between the features learned by a standard SAE (decoder weight vectors) and the data-generating features against the L^1 penalty coefficient in the training loss, following Sharkey et al. (2022). There is a ‘Goldilocks zone’ where SAEs near-perfectly recover the data-generating features, given enough latents to represent them.

to a two-layer MLP at initialization and trained SAEs on both the inputs and outputs individually. As a control, we used Gaussian-distributed inputs with a mean and standard deviation equal to the superposed toy data. We used standard SAEs with an L^1 sparsity penalty (Appendix I.2).

Following Sharkey et al. (2022), we confirmed that SAEs can recover the ground-truth features that generated the data (Figure 4). In particular, we measured the mean max cosine similarity (MMCS): for every data-generating feature, we found its maximum cosine similarity with the features learned by the SAE (its decoder weight vectors) and took the average over data-generating features. However, the MMCS only applies to the MLP inputs, where we have access to the data-generating features – a different approach is required to analyze the MLP outputs. To this end, we took the ability of SAEs to achieve low reconstruction error with high sparsity as a proxy for the degree to which the training data exhibits superposition. Specifically, we vary the L^1 penalty coefficient to obtain Pareto frontiers of the explained variance against sparsity measures (Figure 5a).

As expected, we found that SAEs achieved much greater sparsity at a given level of explained variance for the superposed inputs relative to the Gaussian control (Figure 5a; orange and blue-green). Interestingly, the difference between the superposed outputs, i.e., the outputs of the MLP given the superposed inputs, and the Gaussian outputs is much smaller, with only slightly greater sparsity at a given level of explained variance. This suggests that the outputs of randomly initialized MLPs have a relatively high level of sparsity insensitive to the input distribution. We consider other sparsity measures and hyperparameters in Appendix I.

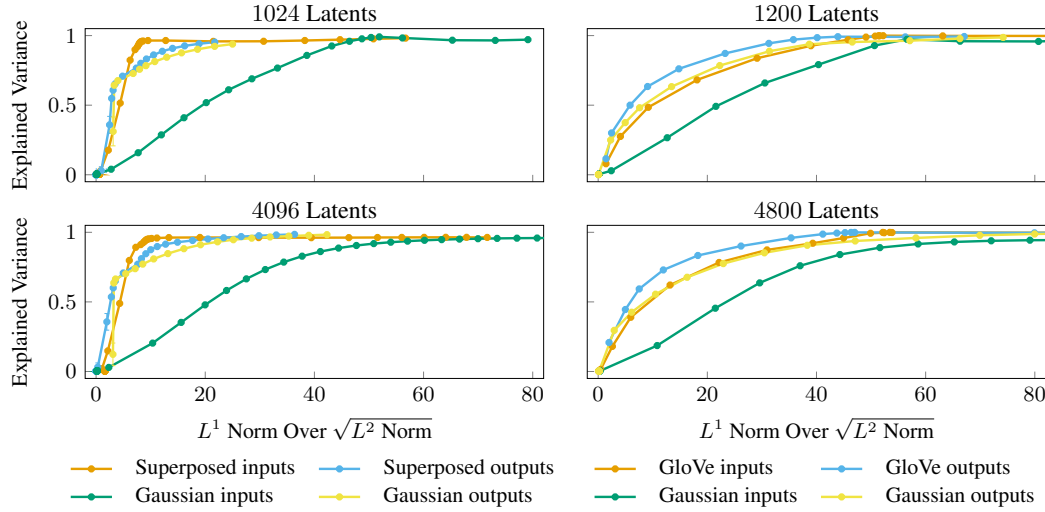
4.3 DO TOKEN EMBEDDINGS EXHIBIT SUPERPOSITION?

To the extent that randomly initialized neural networks preserve or amplify superposition, our results (Section 3) could be explained by the degree to which the inputs to transformer language models exhibit superposition. We study this question by applying the procedure described in Section 4.2 to language data. In particular, we train SAEs on pre-trained GloVe word vectors, the embedding matrices of Pythia models, the results of passing these inputs to a randomly initialized two-layer MLP, and Gaussian controls. The setup is unchanged from Section 4.2, except that the number of data points is fixed by the number of word embeddings or tokens, and we use a single random seed.

We find that the gap between the Pareto frontiers of the GloVe word vectors and the corresponding Gaussian controls (Figure 5b) is smaller than that observed for the toy superposed datasets described in Section 4.2 (Figure 5a). More interestingly, we again see that the Pareto frontiers for both inputs improve when they are passed to a randomly initialized two-layer MLP, emphasizing the possibility that random NNs ‘sparsify’ their inputs (i.e., increase the degree of apparent superposition).

5 LIMITATIONS

In this work, we demonstrate that auto-interpretability measures can produce apparently meaningful, interpretable results for SAEs trained on randomly initialized models, which are unlikely to exhibit



(a) Toy datasets following Sharkey et al. (2022). (b) 300-dim. GloVe vectors (Pennington et al., 2014).

Figure 5: Pareto frontiers of the explained variance against the L^1 norm divided by the square root of the L^2 norm (sparsity) for datasets perhaps exhibiting superposition, Gaussian controls with the same mean and variance, and the corresponding outputs when these are passed to a randomly initialized two-layer MLP (Section 4.2. Each point denotes a choice of L^1 penalty coefficient.

computationally interesting features. Given the impossibility of testing across all datasets and model architectures, we strategically focused on the Pythia family of models, widely adopted in mechanistic interpretability research (e.g., Paulo and Belrose, 2025; Ghilardi et al., 2025; Mueller, 2024), and the RedPajama-V2 dataset, representing typical pre-training data for language models and SAEs.

While we used the default model for generating explanations in the EleutherAI auto-interpretability framework (Caden et al., 2025), exploring alternative models could yield valuable insights into aggregate behaviors and the quality of generated explanations. Importantly, we do not claim that SAEs fail to capture information from trained Transformers above and beyond randomly initialized transformers; only that aggregate auto-interpretability measures do not necessarily indicate the existence of interesting underlying features.

6 CONCLUSION

In this work, we applied sparse autoencoders to both trained and randomly initialized transformers and evaluated them with a suite of common quantitative metrics. Our central empirical finding is that, under certain conditions, these metrics – particularly aggregate auto-interpretability scores – can be surprisingly similar in both settings. While we observe that features derived from trained transformers are qualitatively more complex and abstract, especially in later layers, these aggregate metrics often fail to capture this distinction.

This result does not imply that SAEs trained on real models fail to learn meaningful computational features. Rather, it reveals a limitation in our current evaluation methods. High aggregate auto-interpretability scores are insufficient proof for the discovery of complex, learned computations: they may instead reflect simpler structure inherent in the data or model architecture that is preserved even by random weights. Our analysis of token distribution entropy, while preliminary, serves as a proof-of-concept: it successfully revealed differences in feature ‘abstractness’ that aggregate auto-interpretability scores missed. Future work should focus on developing more robust metrics that can quantify the computational significance of the features SAEs discover. Our work reaffirms the importance of benchmarking interpretability techniques against strong, appropriately constructed null models, such as the randomly initialized transformers used here. Without such baselines, it is difficult to confidently attribute discovered features to the process of learning.

REFERENCES

- J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim. Sanity Checks for Saliency Maps, 2020. URL <https://arxiv.org/abs/1810.03292>.
- A. J. Bell and T. J. Sejnowski. An Information-Maximization Approach to Blind Separation and Blind Deconvolution. *Neural Computation*, 7(6):1129–1159, Nov. 1995. ISSN 0899-7667. doi: 10.1162/neco.1995.7.6.1129. URL <https://ieeexplore.ieee.org/abstract/document/6796129>. Conference Name: Neural Computation.
- N. Belrose, L. Quirke, A. Garriga-Alonso, neverix, HongchuanZeng, A. Duong, lewington, P. Minervini, T. Vogel, T. Fel, Yang, avgalichin, and T. V. Browne. Eleutherai/sparsify. <https://github.com/EleutherAI/sparsify>, sep 22 2025. URL <https://github.com/EleutherAI/sparsify>.
- S. Biderman, H. Schoelkopf, Q. G. Anthony, H. Bradley, K. O’Brien, E. Hallahan, M. A. Khan, S. Purohit, U. S. Prashanth, E. Raff, A. Skowron, L. Sutawika, and O. V. D. Wal. Pythia: A Suite for Analyzing Large Language Models Across Training and Scaling. In *Proceedings of the 40th International Conference on Machine Learning*, pages 2397–2430. PMLR, July 2023. URL <https://proceedings.mlr.press/v202/biderman23a.html>. ISSN: 2640-3498.
- S. Bills, N. Cammarata, D. Mossing, H. Tillman, L. Gao, G. Goh, I. Sutskever, J. Leike, J. Wu, and W. Saunders. Language models can explain neurons in language models, May 2023. URL <https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html>.
- D. Braun, J. Taylor, N. Goldowsky-Dill, and L. Sharkey. Identifying Functionally Important Features with End-to-End Sparse Dictionary Learning, May 2024. URL <http://arxiv.org/abs/2405.12241>. arXiv:2405.12241 [cs].
- T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner, C. Anil, C. Denison, and A. Askell. Towards Monosemanticity: Decomposing Language Models With Dictionary Learning, 2023. URL <https://transformer-circuits.pub/2023/monosemantic-features>.
- Caden, L. Quirke, G. Paulo, A. Mallen, neverix, N. Belrose, johnny, A. Sumer, A. Duong, W. Li, S. Biderman, L. Farnik, I. E. Ashimine, A. Muhamed, D. Ghilardi, F. Xiao, and M. M. Rahman. Eleutherai/delphi. <https://github.com/EleutherAI/delphi>, sep 22 2025. URL <https://github.com/EleutherAI/delphi>.
- L. Chan. Superposition is not “just” neuron polysemanticity. Apr. 2024. URL <https://www.alignmentforum.org/posts/8EyCQKuWo6swZpagS/superposition-is-not-just-neuron-polysemanticity>.
- M. Chaudhary and A. Geiger. Evaluating Open-Source Sparse Autoencoders on Disentangling Factual Knowledge in GPT-2 Small, Sept. 2024. URL <http://arxiv.org/abs/2409.04478>. arXiv:2409.04478 [cs].
- D. Choi, V. Huang, K. Meng, D. D. Johnson, J. Steinhardt, and S. Schwettmann. Scaling Automatic Neuron Description, Oct. 2024. URL <https://transluce.org/neuron-descriptions>.
- H. Cunningham, A. Ewart, L. Riggs, R. Huben, and L. Sharkey. Sparse Autoencoders Find Highly Interpretable Features in Language Models, Oct. 2023. URL <http://arxiv.org/abs/2309.08600>. arXiv:2309.08600 [cs].
- T. Dooms and D. Wilhelm. Tokenized SAEs: Disentangling SAE Reconstructions. June 2024. URL <https://openreview.net/forum?id=5Eas7HCe38>.
- N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds, R. Lasenby, D. Drain, C. Chen, R. Grosse, S. McCandlish, J. Kaplan, D. Amodei, M. Wattenberg, and C. Olah. Toy Models of Superposition, Sept. 2022. URL <http://arxiv.org/abs/2209.10652>. arXiv:2209.10652 [cs].
- J. Engels, I. Liao, E. J. Michaud, W. Gurnee, and M. Tegmark. Not All Language Model Features Are Linear, May 2024. URL <http://arxiv.org/abs/2405.14860>. arXiv:2405.14860 [cs].

-
- L. Farnik, T. Lawson, C. Houghton, and L. Aitchison. Jacobian Sparse Autoencoders: Sparsify Computations, Not Just Activations. In *Forty-second International Conference on Machine Learning*, June 2025. URL <https://openreview.net/forum?id=TPuFRuNano>.
- A. Foote, N. Nanda, E. Kran, I. Konstas, and F. Barez. N2G: A Scalable Approach for Quantifying Interpretable Neuron Representations in Large Language Models, Apr. 2023. URL <http://arxiv.org/abs/2304.12918>. arXiv:2304.12918 [cs].
- L. Gao, T. D. la Tour, H. Tillman, G. Goh, R. Troll, A. Radford, I. Sutskever, J. Leike, and J. Wu. Scaling and evaluating sparse autoencoders, June 2024. URL <http://arxiv.org/abs/2406.04093>. arXiv:2406.04093 [cs].
- D. Ghilardi, F. Belotti, M. Molinari, T. Ma, and M. Palmonari. Group-SAE: Efficient Training of Sparse Autoencoders for Large Language Models via Layer Groups, Sept. 2025.
- W. Gurnee, N. Nanda, M. Pauly, K. Harvey, D. Troitskii, and D. Bertsimas. Finding Neurons in a Haystack: Case Studies with Sparse Probing, June 2023. URL <http://arxiv.org/abs/2305.01610>. arXiv:2305.01610 [cs].
- S. Heimersheim and N. Nanda. How to use and interpret activation patching, Apr. 2024. URL <http://arxiv.org/abs/2404.15255>. arXiv:2404.15255 [cs].
- J. Huang, Z. Wu, C. Potts, M. Geva, and A. Geiger. RAVEL: Evaluating Interpretability Methods on Disentangling Language Model Representations, Aug. 2024. URL <http://arxiv.org/abs/2402.17700>. arXiv:2402.17700 [cs].
- A. Hyvärinen and E. Oja. Independent component analysis: algorithms and applications. *Neural Networks*, 13(4):411–430, June 2000. ISSN 0893-6080. doi: 10.1016/S0893-6080(00)00026-5. URL <https://www.sciencedirect.com/science/article/pii/S0893608000000265>.
- E. T. Jaynes. *Probability Theory: The Logic of Science*. Cambridge University Press, Cambridge, UK, 2003. ISBN 0521592712. The maximum entropy property of the Gaussian distribution is discussed on pages 208 and 365, in chapters 7 and 11.
- A. Karvonen, C. Rager, J. Lin, C. Tigges, J. Bloom, D. Chanin, Y.-T. Lau, E. Farrell, A. Conmy, C. McDougall, K. Ayonrinde, M. Wearden, S. Marks, and N. Nanda. SAE Bench: A Comprehensive Benchmark for Sparse Autoencoders, Dec. 2024a. URL <https://www.neuronpedia.org/sae-bench/info>.
- A. Karvonen, C. Rager, S. Marks, and N. Nanda. Evaluating Sparse Autoencoders on Targeted Concept Erasure Tasks, Nov. 2024b. URL <https://arxiv.org/abs/2411.18895>.
- A. Karvonen, B. Wright, C. Rager, R. Angell, J. Brinkmann, L. Smith, C. M. Verdun, D. Bau, and S. Marks. Measuring Progress in Dictionary Learning for Language Model Interpretability with Board Game Models. *arXiv preprint arXiv:2408.00113*, 2024c. URL <https://arxiv.org/abs/2408.00113>.
- C. Kissane, R. Krzyzanowski, J. I. Bloom, A. Conmy, and N. Nanda. Interpreting Attention Layer Outputs with Sparse Autoencoders. June 2024. URL <https://openreview.net/forum?id=feWUBDwjji>.
- T. Lawson, L. Farnik, C. Houghton, and L. Aitchison. Residual Stream Analysis with Multi-Layer SAEs. In *The Thirteenth International Conference on Learning Representations*, Oct. 2024. URL <https://openreview.net/forum?id=XAjfjizaKs>.
- V. Lecomte, K. Thaman, R. Schaeffer, N. Bashkansky, T. Chow, and S. Koyejo. What Causes Polysemanticity? An Alternative Origin Story of Mixed Selectivity from Incidental Causes, Feb. 2024. URL <http://arxiv.org/abs/2312.03096>. arXiv:2312.03096 [cs].
- H. Lee, A. Battle, R. Raina, and A. Ng. Efficient sparse coding algorithms. In *Advances in Neural Information Processing Systems*, volume 19. MIT Press, 2006. URL https://proceedings.neurips.cc/paper_files/paper/2006/hash/2d71b2ae158c7c5912cc0bbde2bb9d95-Abstract.html.

594 T. Lieberum, S. Rajamanoharan, A. Conmy, L. Smith, N. Sonnerat, V. Varma, J. Kramár, A. Dragan,
595 R. Shah, and N. Nanda. Gemma Scope: Open Sparse Autoencoders Everywhere All At Once on
596 Gemma 2, Aug. 2024. URL <http://arxiv.org/abs/2408.05147>. arXiv:2408.05147 [cs].
597

598 J. Lin and J. Bloom. Announcing Neuronpedia: Platform for accelerating research into Sparse
599 Autoencoders, Mar. 2024. URL [https://www.alignmentforum.org/posts/BaEQoxHhWPrkinm](https://www.alignmentforum.org/posts/BaEQoxHhWPrkinmxd/announcing-neuronpedia-platform-for-accelerating-research)
600 [xd/announcing-neuronpedia-platform-for-accelerating-research](https://www.alignmentforum.org/posts/BaEQoxHhWPrkinmxd/announcing-neuronpedia-platform-for-accelerating-research).

601 J. Lindsey, A. Templeton, J. Marcus, T. Conerly, and J. Batson. Sparse Crosscoders for Cross-Layer
602 Features and Model Diffing, Oct. 2024. URL [https://transformer-circuits.pub/2024/cro](https://transformer-circuits.pub/2024/crosscoders/index.html)
603 [sscoders/index.html](https://transformer-circuits.pub/2024/crosscoders/index.html).
604

605 L. Liu, X. Liu, J. Gao, W. Chen, and J. Han. Understanding the difficulty of training transformers.
606 In B. Webber, T. Cohn, Y. He, and Y. Liu, editors, *Proceedings of the 2020 Conference on*
607 *Empirical Methods in Natural Language Processing (EMNLP)*, pages 5747–5763, Online, Nov.
608 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.463. URL
609 <https://aclanthology.org/2020.emnlp-main.463/>.

610 A. Makelov. Sparse Autoencoders Match Supervised Features for Model Steering on the IOI Task.
611 June 2024. URL <https://openreview.net/forum?id=JdrVuEQih5>.
612

613 A. Makhzani and B. Frey. k-Sparse Autoencoders, Mar. 2014. URL [http://arxiv.org/abs/1312](http://arxiv.org/abs/1312.5663)
614 [.5663](http://arxiv.org/abs/1312.5663). arXiv:1312.5663 [cs].

615 S. C. Marshall and J. H. Kirchner. Understanding polysemanticity in neural networks through coding
616 theory, Jan. 2024. URL <http://arxiv.org/abs/2401.17975>. arXiv:2401.17975 [cs].
617

618 K. Meng, D. Bau, A. Andonian, and Y. Belinkov. Locating and Editing Factual Associations in
619 GPT. *Advances in Neural Information Processing Systems*, 35:17359–17372, Dec. 2022. URL
620 [https://proceedings.neurips.cc/paper_files/paper/2022/hash/6f1d43d5a82a37e89](https://proceedings.neurips.cc/paper_files/paper/2022/hash/6f1d43d5a82a37e89b0665b33bf3a182-Abstract-Conference.html)
621 [b0665b33bf3a182-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/6f1d43d5a82a37e89b0665b33bf3a182-Abstract-Conference.html).

622 A. Mudide, J. Engels, E. J. Michaud, M. Tegmark, and C. S. de Witt. Efficient Dictionary Learning
623 with Switch Sparse Autoencoders. *arXiv preprint arXiv:2410.08201*, 2024. URL [https://arxiv.](https://arxiv.org/abs/2410.08201)
624 [org/abs/2410.08201](https://arxiv.org/abs/2410.08201).
625

626 A. Mueller. Missed Causes and Ambiguous Effects: Counterfactuals Pose Challenges for Interpreting
627 Neural Networks, 2024. URL <https://arxiv.org/abs/2407.04690>.

628 A. Ng. Sparse autoencoder, 2011. URL [https://graphics.stanford.edu/courses/cs233-2](https://graphics.stanford.edu/courses/cs233-21-spring/ReferencedPapers/SAE.pdf)
629 [1-spring/ReferencedPapers/SAE.pdf](https://graphics.stanford.edu/courses/cs233-21-spring/ReferencedPapers/SAE.pdf).
630

631 K. O’Brien, D. Majercak, X. Fernandes, R. Edgar, J. Chen, H. Nori, D. Carignan, E. Horvitz, and
632 F. Poursabzi-Sangde. Steering Language Model Refusal with Sparse Autoencoders, Nov. 2024.
633 URL <https://arxiv.org/abs/2411.11296>.

634 B. A. Olshausen and D. J. Field. Emergence of simple-cell receptive field properties by learning a
635 sparse code for natural images. *Nature*, 381(6583):607–609, June 1996. ISSN 1476-4687. doi:
636 10.1038/381607a0. URL <https://www.nature.com/articles/381607a0>. Publisher: Nature
637 Publishing Group.
638

639 B. A. Olshausen and D. J. Field. Sparse coding with an overcomplete basis set: A strategy employed
640 by V1? *Vision Research*, 37(23):3311–3325, Dec. 1997. ISSN 0042-6989. doi: 10.1016/S0042-6
641 989(97)00169-7. URL [https://www.sciencedirect.com/science/article/pii/S0042698](https://www.sciencedirect.com/science/article/pii/S0042698997001697)
642 [997001697](https://www.sciencedirect.com/science/article/pii/S0042698997001697).

643 K. Park, Y. J. Choe, and V. Veitch. The Linear Representation Hypothesis and the Geometry of Large
644 Language Models, Nov. 2023. URL <http://arxiv.org/abs/2311.03658>. arXiv:2311.03658
645 [cs, stat].
646

647 G. Paulo and N. Belrose. Sparse Autoencoders Trained on the Same Data Learn Different Features,
2025. URL <https://arxiv.org/abs/2501.16615>.

-
- G. Paulo, A. Mallen, C. Juang, and N. Belrose. Automatically Interpreting Millions of Features in Large Language Models, Oct. 2024.
- J. Pennington, R. Socher, and C. D. Manning. GloVe: Global Vectors for Word Representation. In A. Moschitti, B. Pang, and W. Daelemans, editors, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, Oct. 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1162.
- S. Rajamanoharan, A. Conmy, L. Smith, T. Lieberum, V. Varma, J. Kramar, R. Shah, and N. Nanda. Improving Sparse Decomposition of Language Model Activations with Gated Sparse Autoencoders. June 2024a. URL <https://openreview.net/forum?id=Ppj5KvzU8Q>.
- S. Rajamanoharan, T. Lieberum, N. Sonnerat, A. Conmy, V. Varma, J. Kramár, and N. Nanda. Jumping Ahead: Improving Reconstruction Fidelity with JumpReLU Sparse Autoencoders, July 2024b. URL <http://arxiv.org/abs/2407.14435>. arXiv:2407.14435 [cs].
- A. Scherlis, K. Sachan, A. S. Jermyn, J. Benton, and B. Shlegeris. Polysemanticity and Capacity in Neural Networks, July 2023. URL <http://arxiv.org/abs/2210.01892>. arXiv:2210.01892 [cs].
- L. Sharkey, D. Braun, and B. Millidge. Taking features out of superposition with sparse autoencoders, Dec. 2022. URL <https://www.alignmentforum.org/posts/z6QQJbtpkEAX3AoJJ/interim-research-report-taking-features-out-of-superposition>.
- A. Templeton, T. Conerly, J. Marcus, J. Lindsey, T. Bricken, B. Chen, A. Pearce, C. Citro, E. Ameisen, A. Jones, H. Cunningham, N. L. Turner, C. McDougall, M. MacDiarmid, A. Tamkin, E. Durmus, T. Hume, F. Mosconi, C. D. Freeman, T. R. Sumers, E. Rees, J. Batson, A. Jermyn, S. Carter, C. Olah, and T. Henighan. Scaling Monosemanticity: Extracting Interpretable Features from Claude 3 Sonnet, May 2024. URL <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>.
- M. Wattenberg and F. Viégas. Relational Composition in Neural Networks: A Survey and Call to Action. June 2024. URL <https://openreview.net/forum?id=zzCEiUIPk9>.
- M. Weber, D. Fu, Q. Anthony, Y. Oren, S. Adams, A. Alexandrov, X. Lyu, H. Nguyen, X. Yao, V. Adams, B. Athiwaratkun, R. Chalamala, K. Chen, M. Ryabinin, T. Dao, P. Liang, C. Ré, I. Rish, and C. Zhang. RedPajama: an Open Dataset for Training Large Language Models, 2024. URL <https://arxiv.org/abs/2411.12372>.
- Z. Yun, Y. Chen, B. Olshausen, and Y. LeCun. Transformer visualization via dictionary learning: contextualized embedding as a linear superposition of transformer factors. In E. Agirre, M. Apidianaki, and I. Vulić, editors, *Proceedings of Deep Learning Inside Out (DeeLIO): The 2nd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pages 1–10, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.deelio-1.1. URL <https://aclanthology.org/2021.deelio-1.1>.
- F. Zhang and N. Nanda. Towards Best Practices of Activation Patching in Language Models: Metrics and Methods. In *The Twelfth International Conference on Learning Representations*, Oct. 2023. URL <https://openreview.net/forum?id=Hf17y6u9BC>.
- Z. Zhong and J. Andreas. Algorithmic Capabilities of Random Transformers, 2024. URL <https://arxiv.org/abs/2410.04368>.

A BROADER IMPACT

This work investigates a method currently used for mechanistic interpretability of LLMs, yielding results that challenge certain assumptions about sparse autoencoders. By demonstrating that SAEs can produce similar aggregate auto-interpretability scores for both random and trained transformers, our findings raise important questions about what these SAE evaluation methods are actually capturing.

By better understanding the metrics of SAE quality, we hope that this work will contribute to a more informed search of better SAE-like methods and thus help to make these models more interpretable and to mitigate the potential harm these models could cause. Since our work is an empirical study of the capabilities of a presently used method, and it shows that the method provides interpretation of both random and trained transformers, we think the risk that this work could lead to negative social impact is minimal.

B AUTO-INTERPRETABILITY ROC CURVES

Figures 6, 8, 12 show the similarity between ‘fuzzing’ AUROC for the trained and randomized SAEs for the 70M, 160M, and 1B models. Figures 7, 9, 13, show the similarity between ‘detection’ AUROC for the trained and randomized SAEs for the 70M, 160M, and 1B models.

B.1 PYTHIA 70M

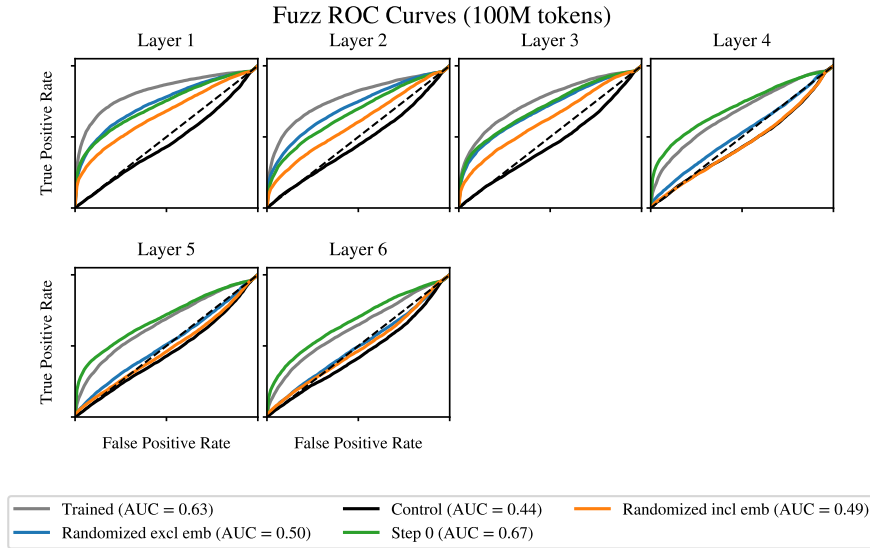


Figure 6: ROC curves for ‘fuzzing’ auto-interpretability for Pythia-70m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases.

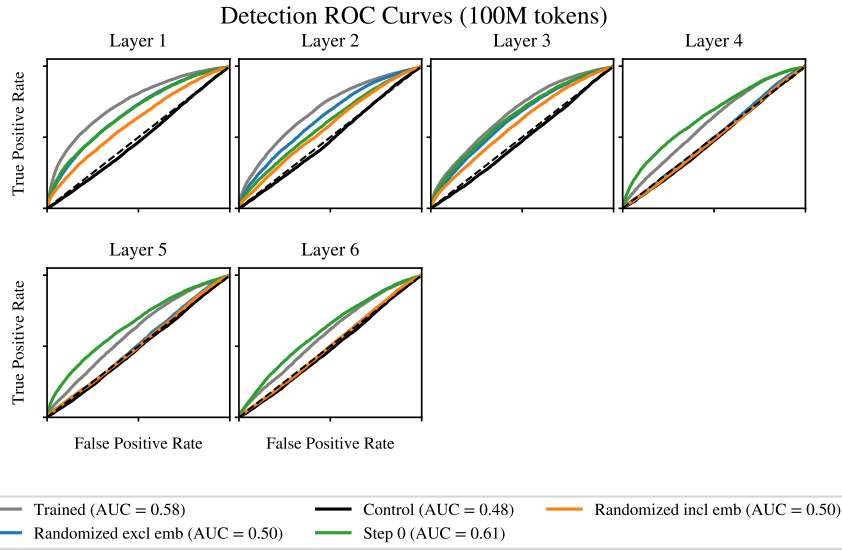


Figure 7: ROC curves for ‘detection’ auto-interpretability for Pythia-70m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases.

B.2 PYTHIA 160M

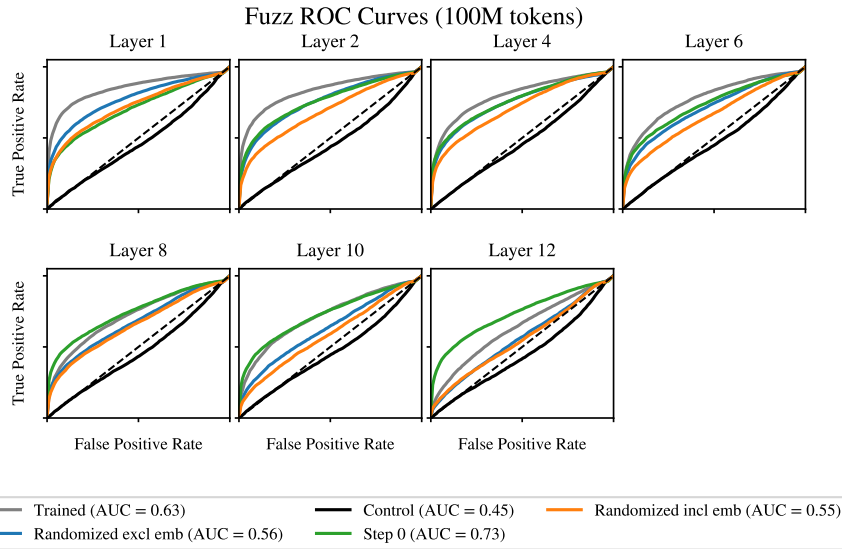


Figure 8: ROC curves for ‘fuzzing’ auto-interpretability for Pythia-160m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases.

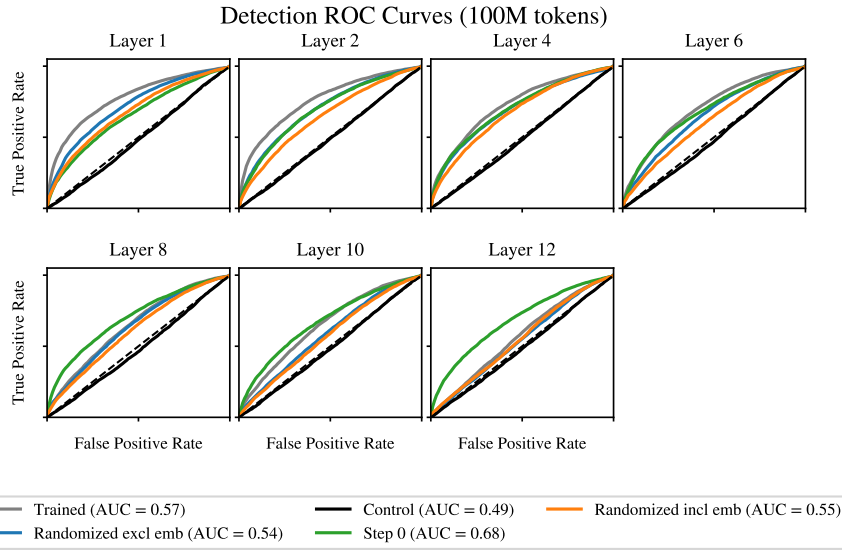


Figure 9: ROC curves for ‘detection’ auto-interpretability for Pythia-160m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases.

B.3 PYTHIA 410M

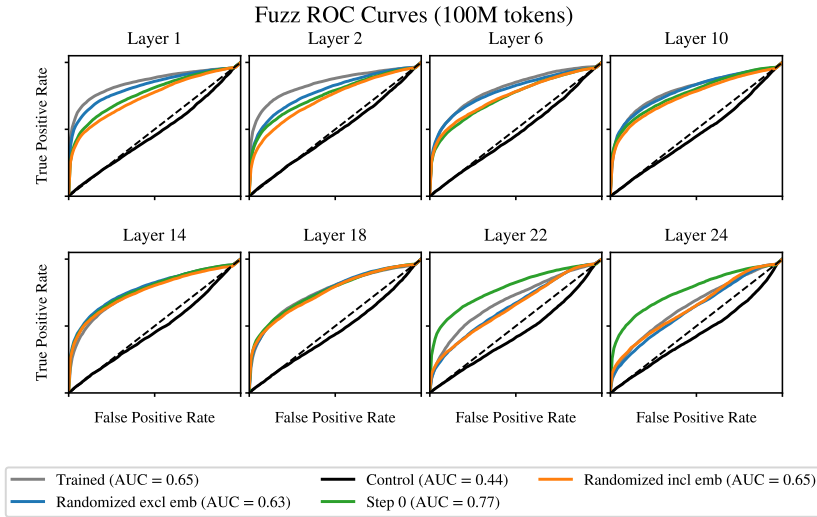


Figure 10: ROC curves for ‘fuzzing’ auto-interpretability for Pythia-410m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases. The auto-interpretability scores here fail to distinguish between trained and randomized models.

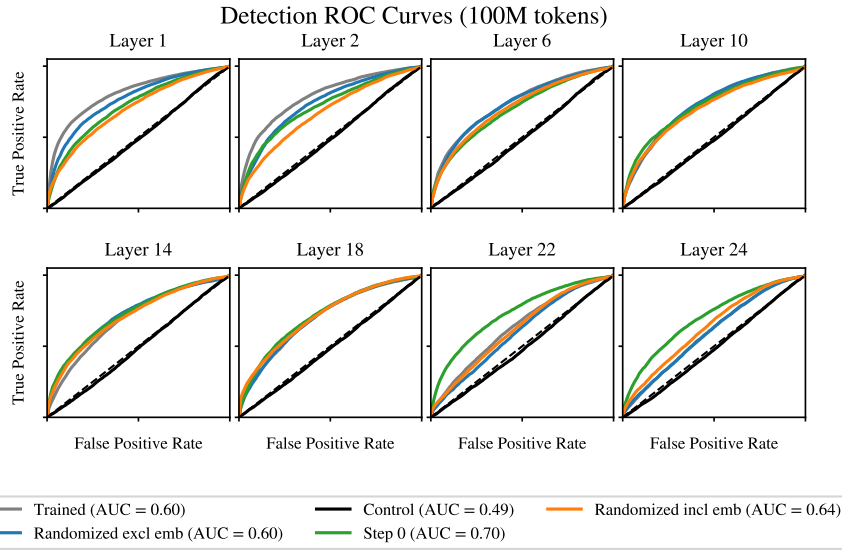


Figure 11: ROC curves for ‘detection’ auto-interpretability for Pythia-410m over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, as well as the overall degradation in performance as the layer index increases.

B.4 PYTHIA-1B

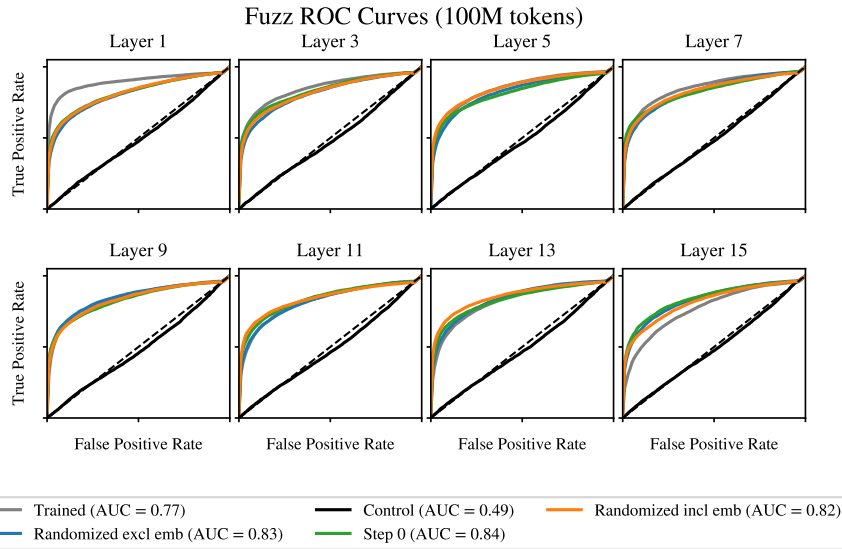


Figure 12: ROC curves for ‘fuzzing’ auto-interpretability for Pythia-1b over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, although here we do not observe an overall degradation in quality.

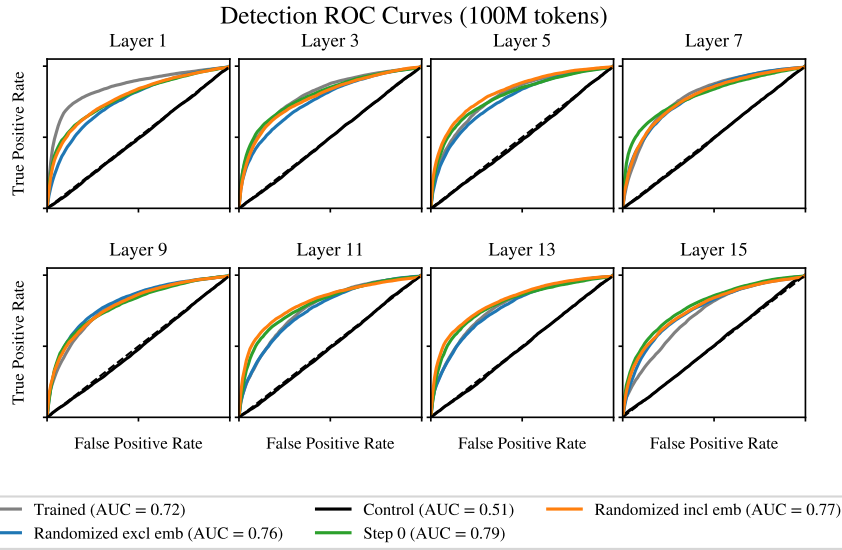


Figure 13: ROC curves for ‘detection’ auto-interpretability for Pythia-1b over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants, although here we do not observe an overall degradation in quality.

B.5 PYTHIA 6.9B

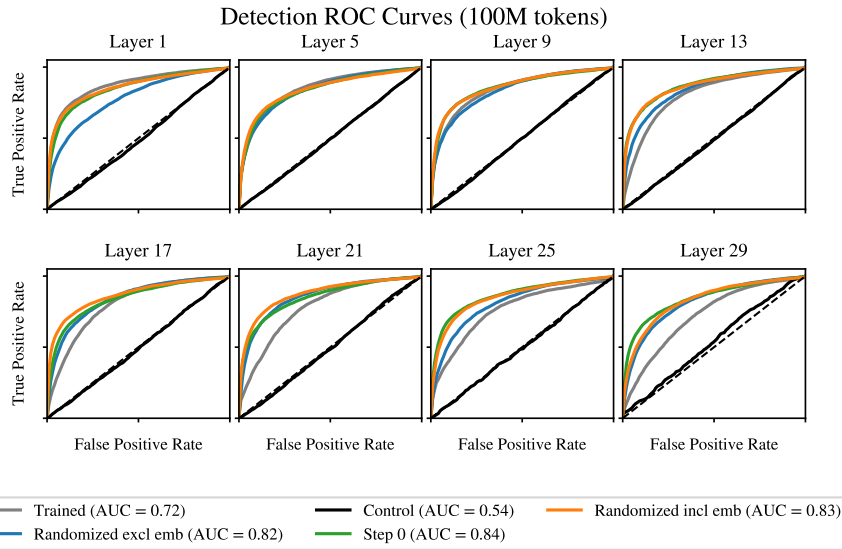


Figure 14: ROC curves for ‘detection’ auto-interpretability for Pythia-6.9b over 100 SAE latents. These results demonstrate the similarity in performance between the SAE variants.

C EFFECT OF INCREASED TRAINING DATA

For our primary experiments, we trained SAEs on 100M tokens (Section 3). We verified that our results were not explained by a lack of sufficient training data by repeating a subset of these experiments with SAEs trained on 1B tokens from the RedPajama dataset (Figure 15).

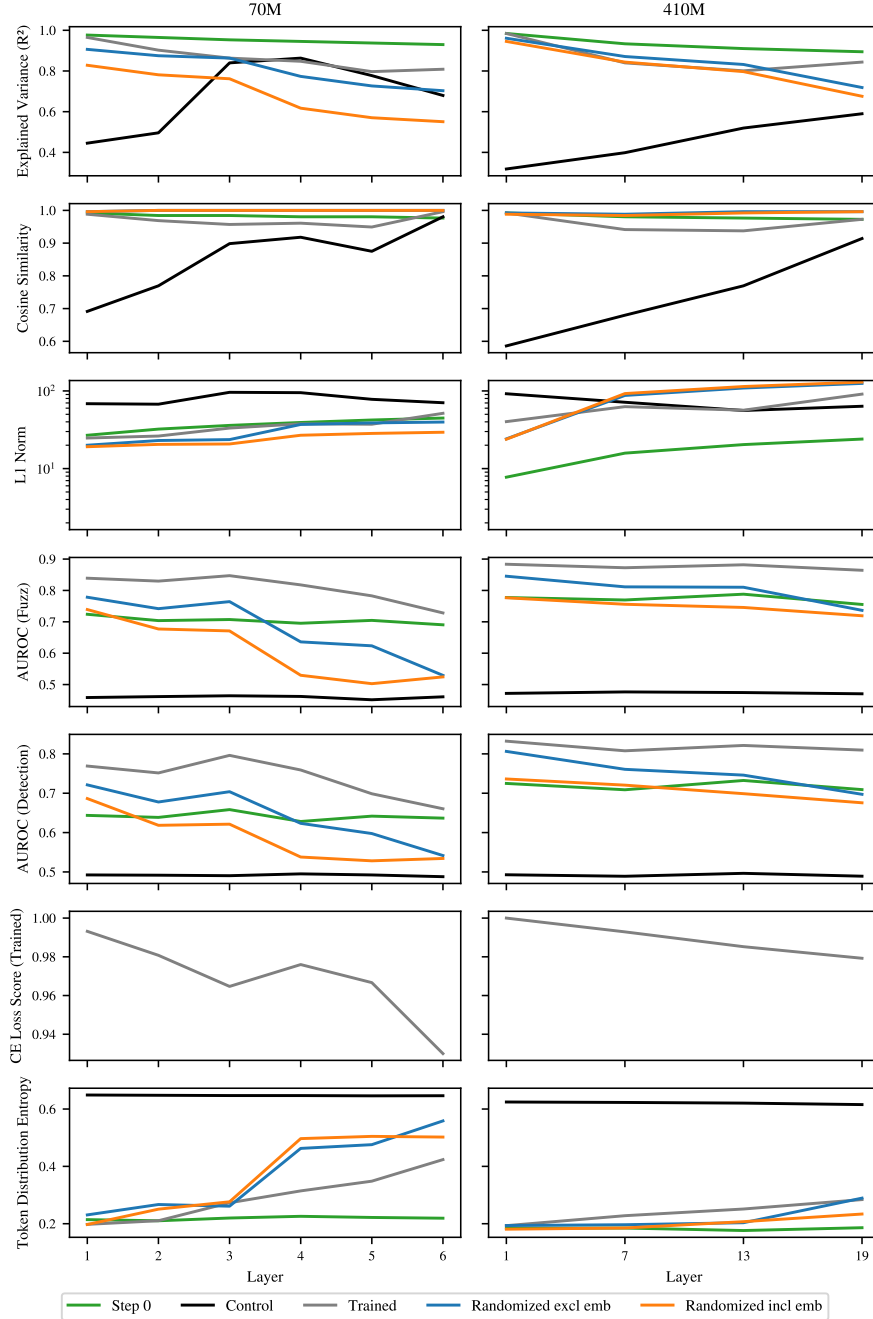
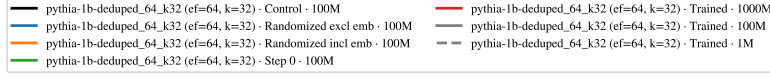


Figure 15: Evaluation metrics for SAEs trained with one billion tokens on the Pythia-70m and 410m models. These results correspond to columns of Figure 2, which show the same evaluation metrics for SAEs trained on 100M tokens, and qualitatively similar behavior.

D EFFECT OF DECREASED TRAINING DATA FOR PYTHIA-1B



Pythia 1B Metrics & ROC Comparison - Token Count Variants (1000M, 100M, 1M)

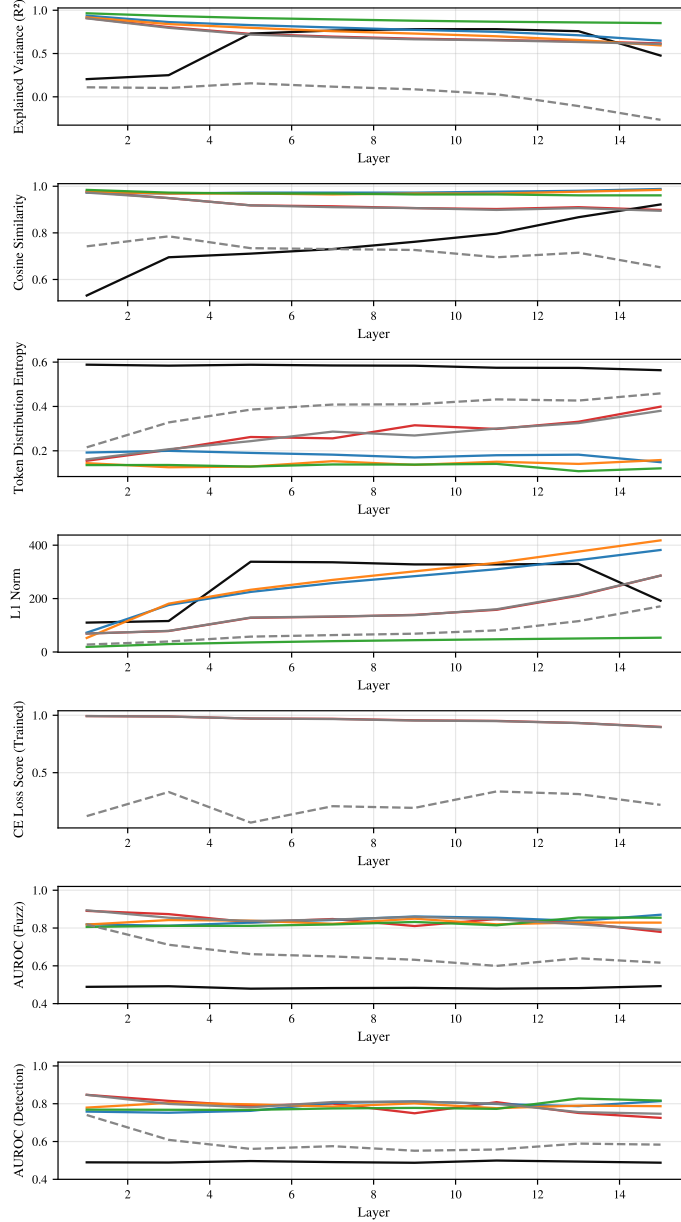


Figure 16: Evaluation metrics for SAEs trained with 1M and 1B tokens on Pythia-1b. The explained variance and CE loss score are significantly lower for the 1M model, showing that the SAEs are under-trained. Average auto-interpretability scores are slightly lower for the earliest layers, but decline sharply with increasing layer. The trends in auto-interpretability and token distribution entropy with layer index are consistent with other SAEs.

E UNCERTAINTY PLOTS FOR PYTHIA-70M

We computed uncertainty for our evaluation metrics on Pythia-70m using five random seeds.

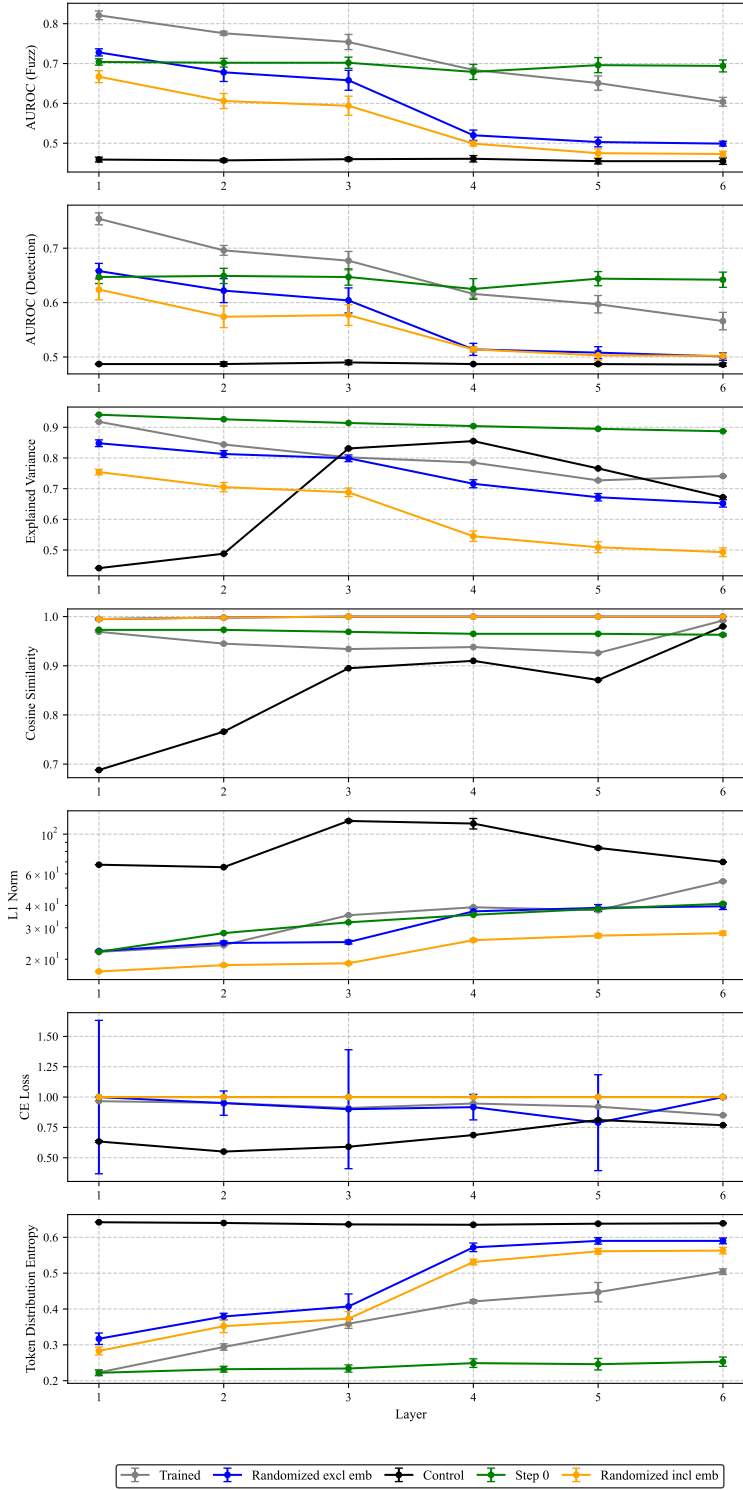


Figure 17: Uncertainty for Pythia-70m metrics computed using five random seeds.

F EFFECT OF SAE HYPERPARAMETERS FOR PYTHIA-160M

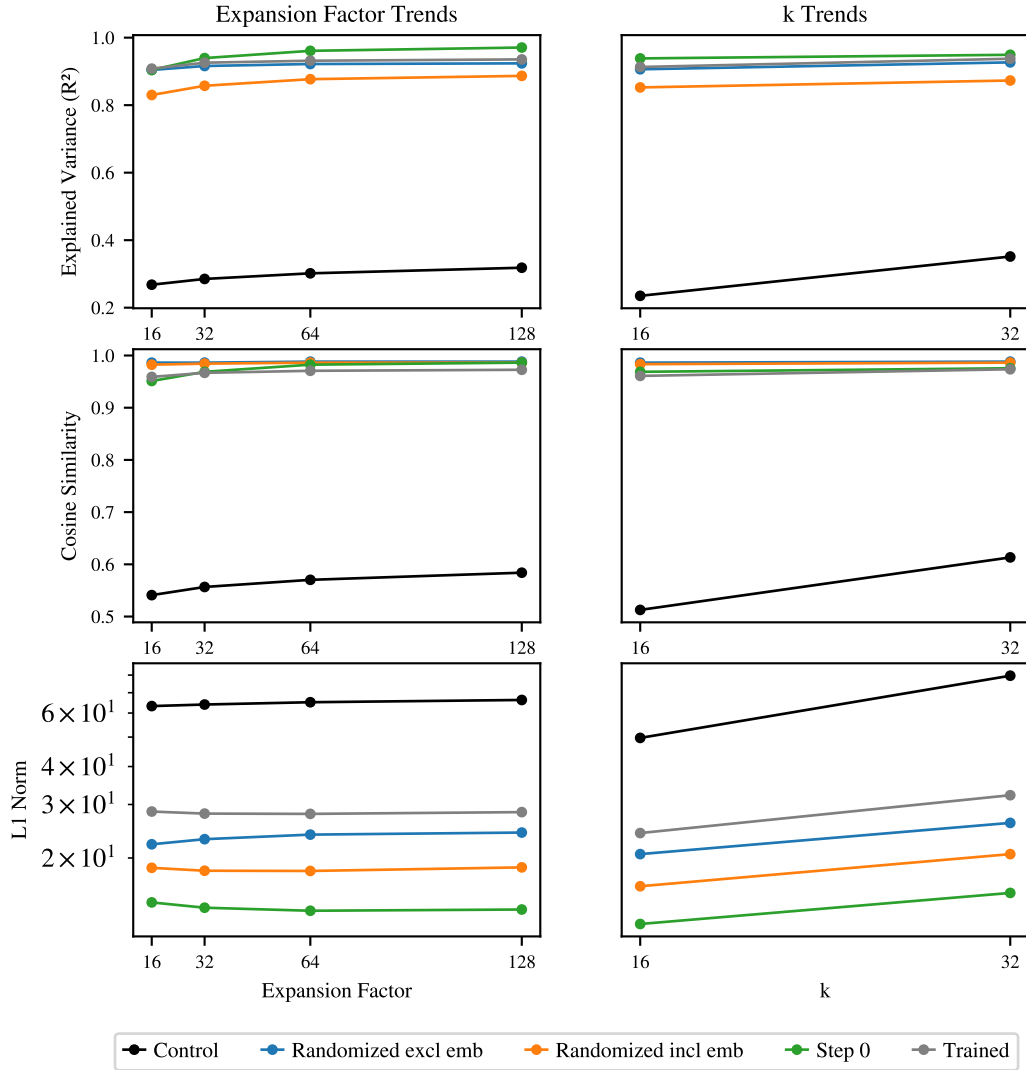


Figure 18: Robustness of SAE performance to hyperparameter selection. Standard evaluation metrics remain stable across a wide range of expansion factors R (16 to 128) and sparsities k (16 to 32), with all initialization strategies maintaining their relative performance ordering. This stability suggests that moderate hyperparameter values (e.g., expansion factor $R = 64$, sparsity $k = 32$) suffice.

G EFFECT OF SAE HYPERPARAMETERS FOR PYTHIA-1B

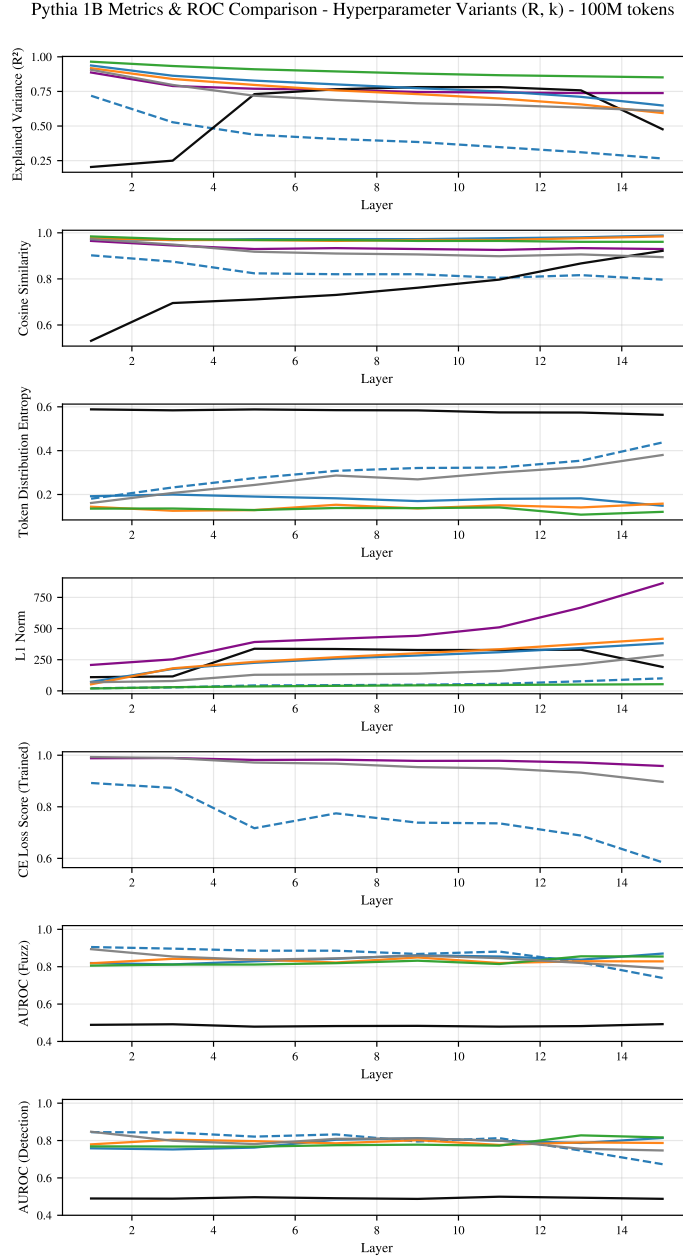


Figure 19: Evaluation metrics for SAEs trained on the Pythia-1b model with different hyperparameters, including the main results from Figure 2. SAEs with a very small expansion factor $R = 2$ and sparsity $k = 4$ are clearly distinguished from our default hyperparameters by the explained variance and CE loss score. Importantly, the auto-interpretability scores of these SAEs remain similar to those trained with default hyperparameters on either trained scores or randomised models.

H TOKEN DISTRIBUTION ENTROPY VS. AUTO-INTERPRETABILITY

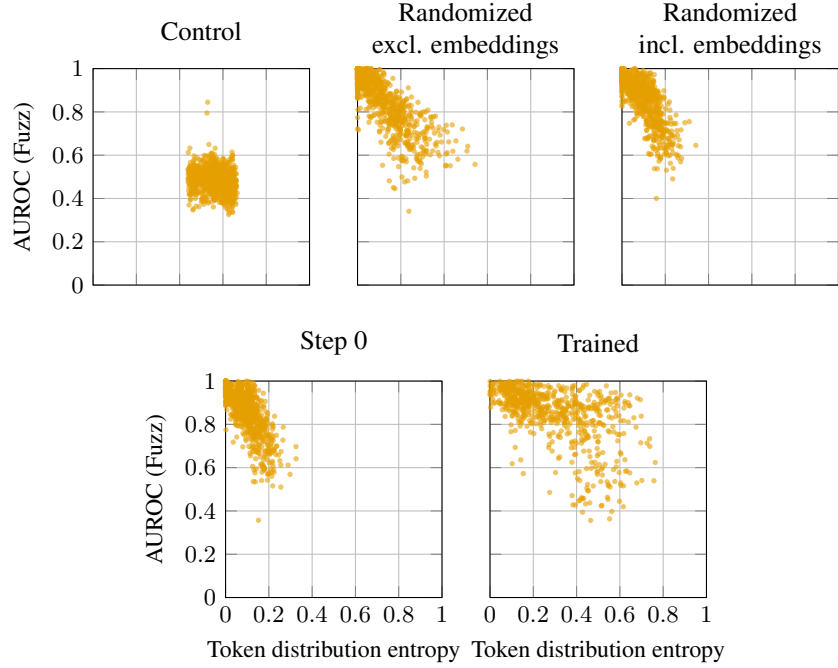


Figure 20: Scatter plots of the per-latent token distribution entropy against ‘fuzzing’ AUROC (auto-interpretability score) for SAEs trained on multiple layers of the Pythia-6.9b model. Each point corresponds to a single latent, taken from the sample of latents used to compute the aggregate metrics displayed in Figures 1 and 2.

Figure 20 clearly distinguishes the negative control, randomized variants, and the trained variant:

- **Control:** Latents have a consistently high entropy (i.e., max activating examples with activation patterns spread across many tokens) and low auto-interpretability score (i.e., generated explanations that fail to adequately explain these activation patterns). No correlation between the two variables is evident.
- **Randomized:** For each of the randomization schemes described in Section 3, we see a negative correlation between entropy and auto-interpretability: in general, the wider variety of tokens for which a latent is activated, the less well the latent’s activation patterns are explained by its generated explanation.
- **Trained:** There is a weaker correlation between the two variables. Crucially, in addition to the broad trend observed for the randomized variants, we also see latents with high entropy *and* auto-interpretability. Some latents have activation patterns that are spread across multiple tokens, which are nevertheless consistent with the latent’s generated explanation.

These results are consistent with the view that aggregate auto-interpretability scores obscure the differences between SAEs based on trained and randomized models. While randomized models with consistent token embeddings can produce ‘single-token’ features, whose activation patterns are easy to explain, only Transformers trained on natural language produce more complex semantic features.

I A TOY MODEL OF SUPERPOSITION

In Section 4, we trained SAEs on toy data designed to exhibit superposition (Sharkey et al., 2022) and GloVe word vectors (Pennington et al., 2014). In this section, we detail the data-generation procedure and training setup.

I.1 DATA GENERATION

First, we construct ground-truth features by sampling n_s points on an n_d -dimensional hypersphere.

For each sample, we determine the feature coefficients by generating $A \in \mathbb{R}^{n_s \times n_s}$ where $A_{ij} \sim \mathcal{N}(0, 1)$, defining a covariance matrix $\Sigma = AA^T$, sampling $\vec{\alpha} \in \mathbb{R}^{n_s}$ where $\alpha_i \sim \mathcal{N}(\vec{0}, \Sigma)$, projecting α_i onto the c.d.f. of $\mathcal{N}(0, 1)$, decaying $\alpha_i \rightarrow \alpha_i^{\lambda_i}$ where $\lambda \in \mathbb{R}$, normalizing $\alpha_i \rightarrow m\alpha_i / n_s \sum_j \alpha_j$ where $m \in \mathbb{R}$, and performing n_s independent Bernoulli trials with $p = \alpha_i$. Finally, we multiply the trial outcomes by n_s independent samples from a continuous uniform distribution $\mathcal{U}_{[0,1)}$.

The parameter λ determines how sharply the frequency of nonzero ground-truth feature coefficients decays with the feature index i . The parameter m is the expected value of the number of nonzero feature coefficients for each sample.

Like Sharkey et al. (2022), we choose $n_s = 512$, $n_d = 256$, $\lambda = 0.99$, and $m = 5$. We include a Python implementation of this procedure in Figure 21.

I.2 TRAINING

The SAEs described in Section 4 comprise a linear encoder with a bias term, a ReLU activation function, and a linear decoder without a bias term. We use orthogonal initialization for the decoder weights and normalize the decoder weight vectors before each training step.

The training loss is the mean squared error (MSE) between the input and decoded vectors, plus the mean L^1 norm of the encoded vectors multiplied by a coefficient, which we vary between 1×10^{-3} and 100.

For the toy data, we train for 100 epochs on 10K data points with 10 random seeds. For the word vectors, we train for 100 epochs on 400K data points with 1 random seed. In both experiments, we reserve 10% of the data points as a validation set, which we use to compute evaluation metrics.

The MLPs described in Section 4 comprise two layers (i.e., one hidden layer) and a ReLU activation function. The input and output sizes are both equal to n_d , and the hidden size is $4n_d$. We loosely based these choices on the feed-forward network components of transformer language models.

```

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362 1 def generate_sharkey(
1363 2     num_samples: int,
1364 3     num_inputs: int,
1365 4     num_features: int,
1366 5     avg_active_features: float,
1367 6     lambda_decay: float,
1368 7 ) -> tuple[Tensor, Tensor]:
1369 8     """
1370 9     Args:
1371 10         num_samples (int): The number of samples to generate.
1372 11         num_inputs (int): The number of input dimensions.
1373 12         num_features (int): The number of ground truth features.
1374 13         avg_active_features (float): The average number of
1375 14             ground truth features active at a time.
1376 15         lambda_decay (float): The exponential decay factor for
1377 16             feature probabilities.
1378 17     """
1379 18     features = torch.randn(num_inputs, num_features)
1380 19     features /= torch.norm(features, dim=0, keepdim=True)
1381 20
1382 21     covariance = torch.randn(num_features, num_features)
1383 22     covariance = covariance @ covariance.T
1384 23     correlated_normal = MultivariateNormal(
1385 24         torch.zeros(num_features), covariance_matrix=covariance
1386 25     )
1387 26
1388 27     samples = []
1389 28     for _ in range(num_samples):
1390 29         p = STANDARD_NORMAL.cdf(correlated_normal.sample())
1391 30         p = p ** (lambda_decay * torch.arange(num_features))
1392 31         p = p * (avg_active_features / (num_features * p.mean()))
1393 32         p = torch.bernoulli(p.clamp(0, 1))
1394 33         coef = p * torch.rand(num_features)
1395 34
1396 35         sample = coef @ features.T
1397 36         samples.append(sample)
1398 37
1399 38     return torch.stack(samples), features
1400
1401
1402
1403

```

Figure 21: A Python implementation of the data-generation procedure introduced by Sharkey et al. (2022) and used in Section 4.

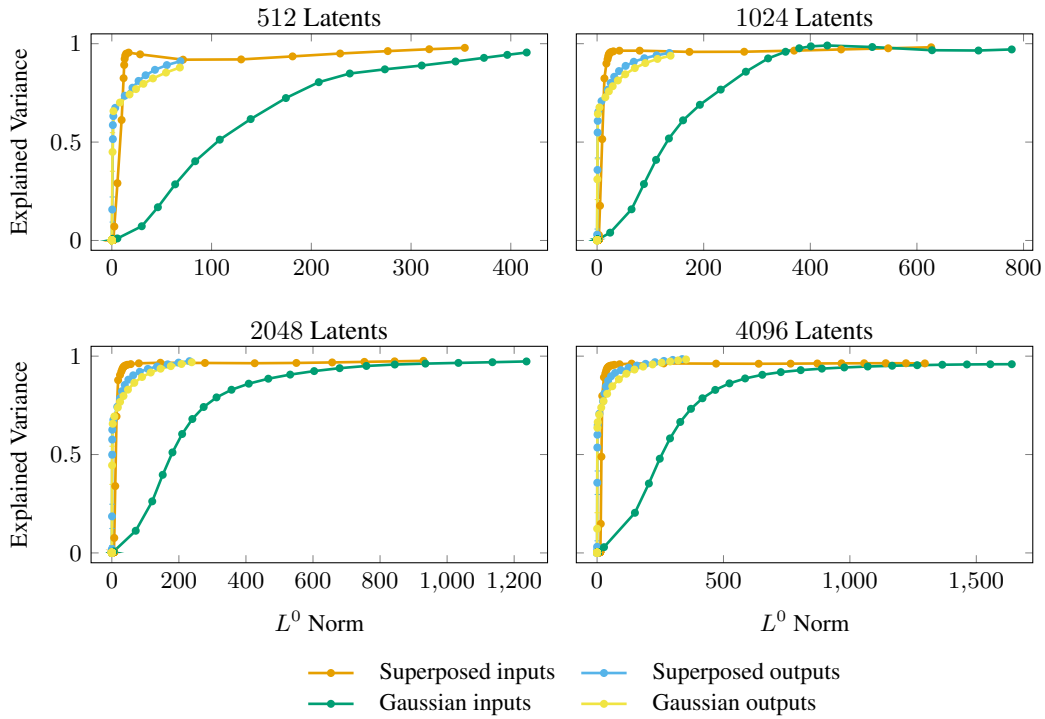


Figure 22: Pareto frontiers of the explained variance against the L^0 norm (sparsity) for toy datasets generated to exhibit superposition, Gaussian controls with the same mean and variance, and the corresponding outputs when these are passed to a randomly initialized two-layer MLP.

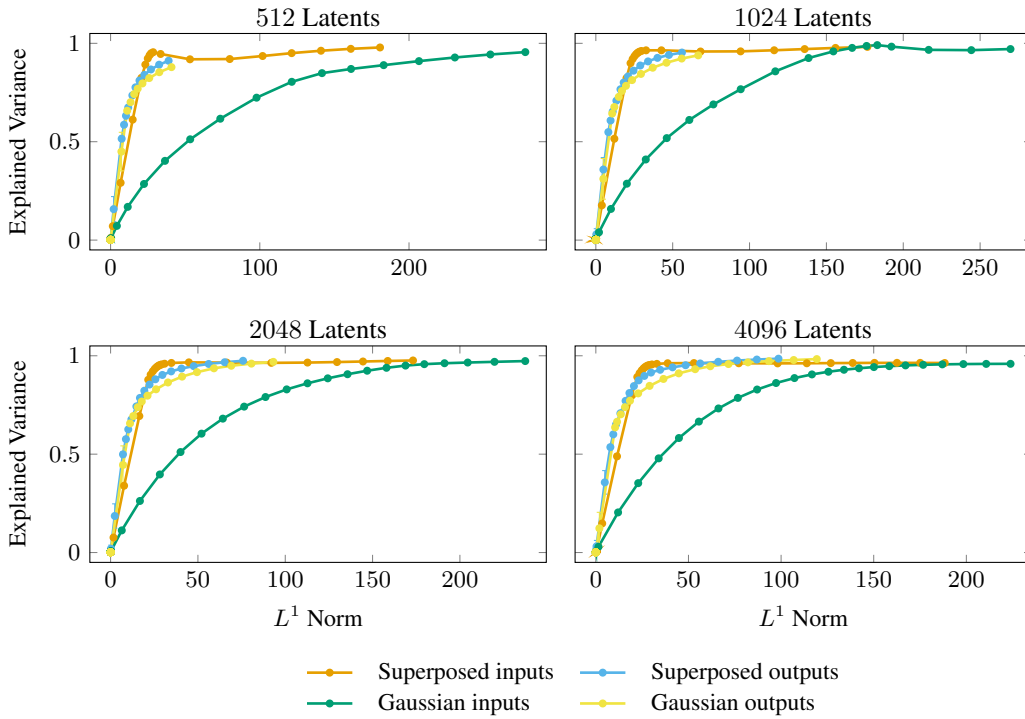


Figure 23: Pareto frontiers of the explained variance against the L^1 norm (sparsity) for toy datasets generated to exhibit superposition, Gaussian controls with the same mean and variance, and the corresponding outputs when these are passed to a randomly initialized two-layer MLP.

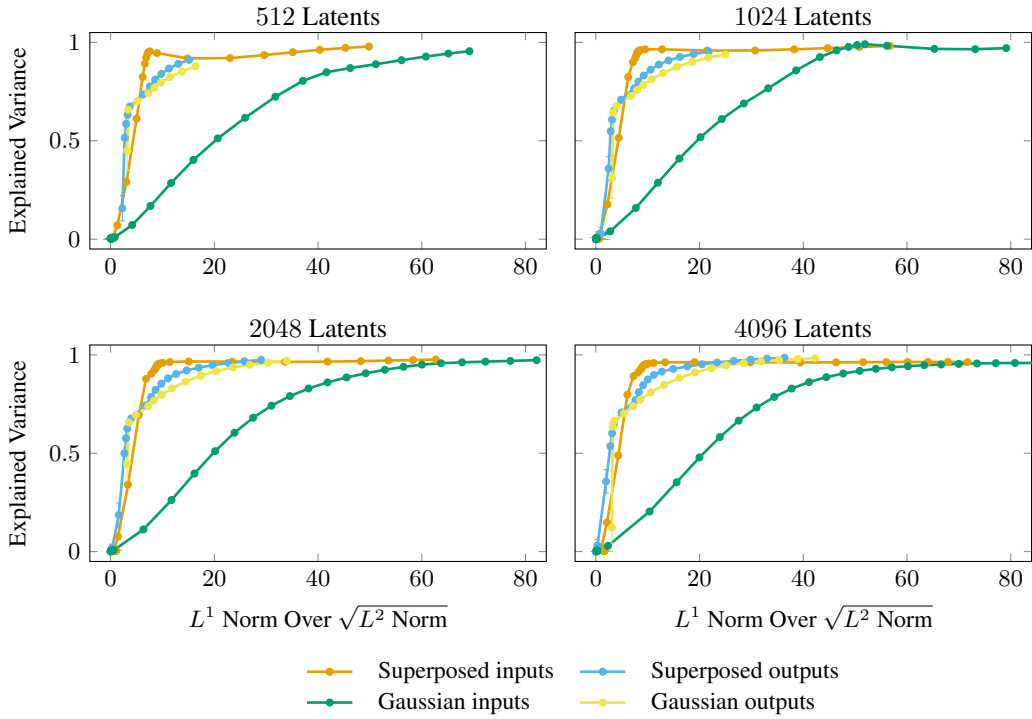


Figure 24: Pareto frontiers of explained variance against the L^1 norm over the square root of the L^2 norm (sparsity) for toy datasets generated to exhibit superposition, Gaussian controls with the same mean and variance, and the corresponding outputs when these are passed to a randomly initialized two-layer MLP.

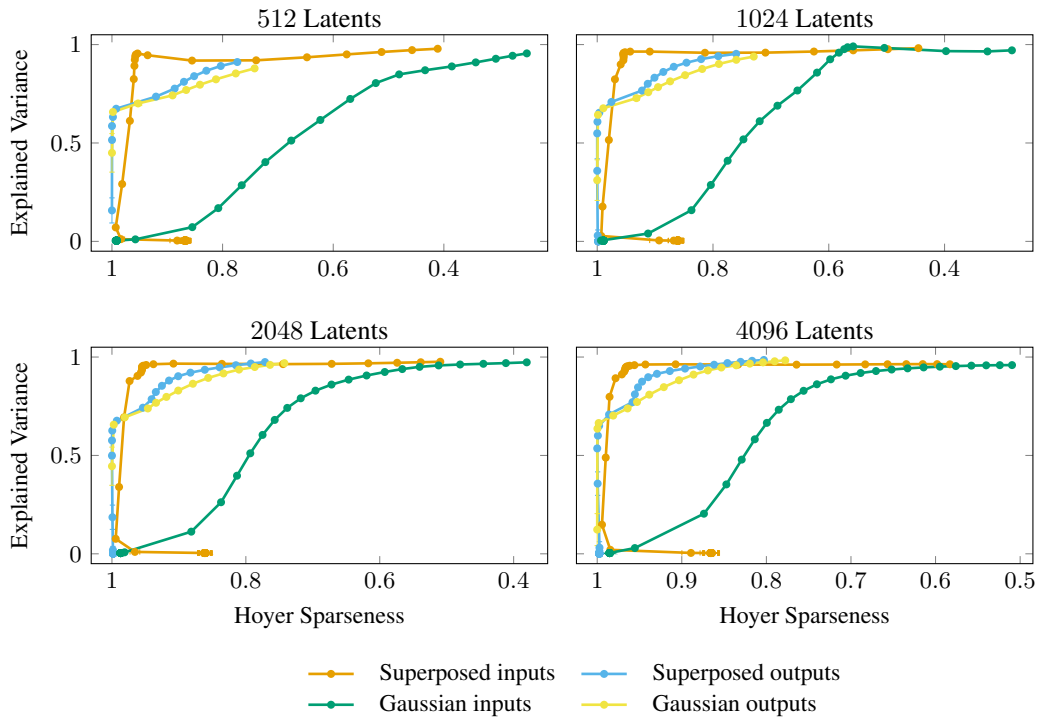


Figure 25: Pareto frontiers of explained variance against the Hoyer sparseness (sparsity) for toy datasets generated to exhibit superposition, Gaussian controls with the same mean and variance, and the corresponding outputs when these are passed to a randomly initialized two-layer MLP.

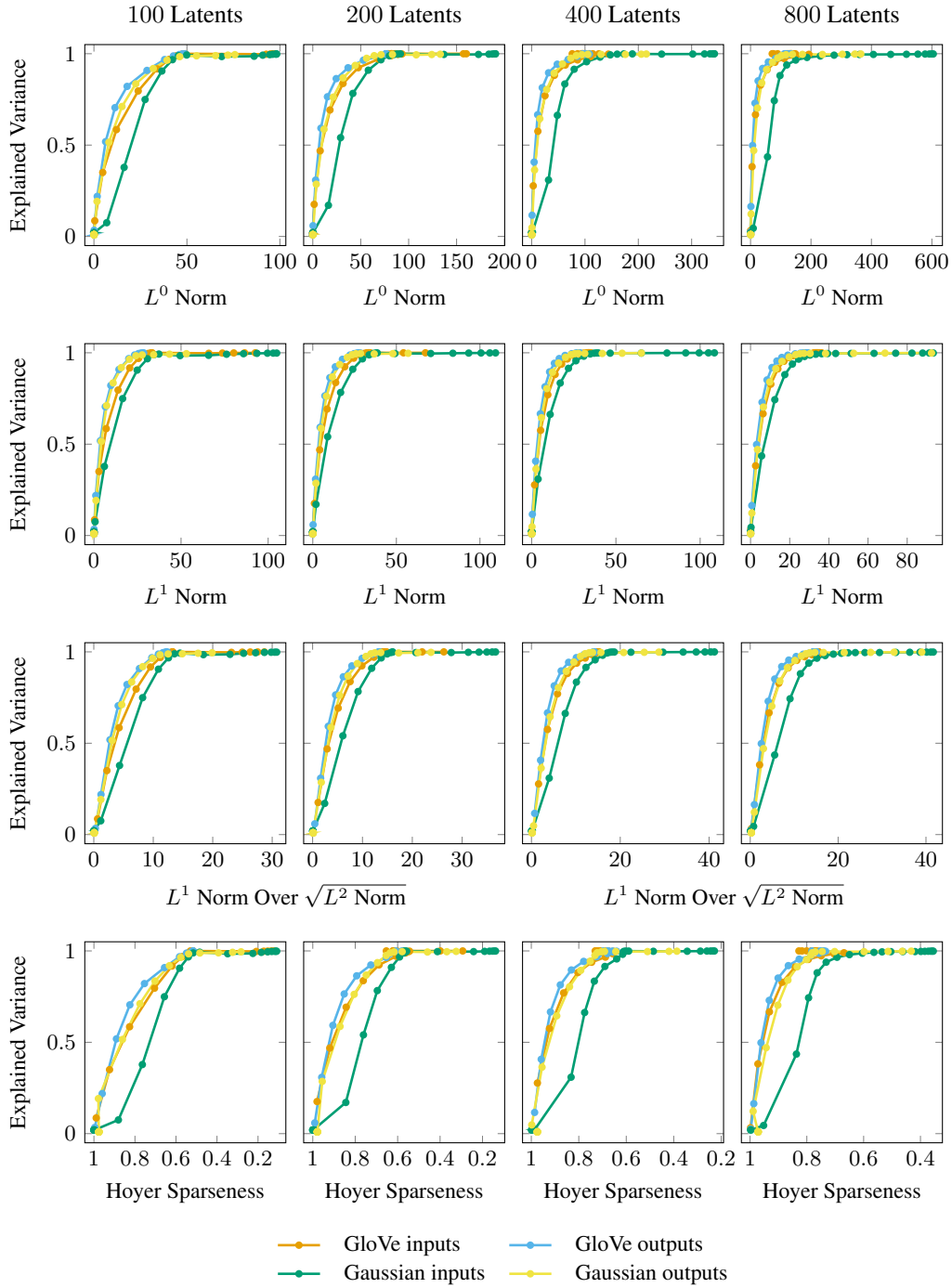


Figure 26: Pareto frontiers of explained variance against sparsity measures for 50-dimensional GloVe word vectors, Gaussian controls with the same mean and variance, and the corresponding outputs when these are passed to a randomly initialized two-layer MLP.

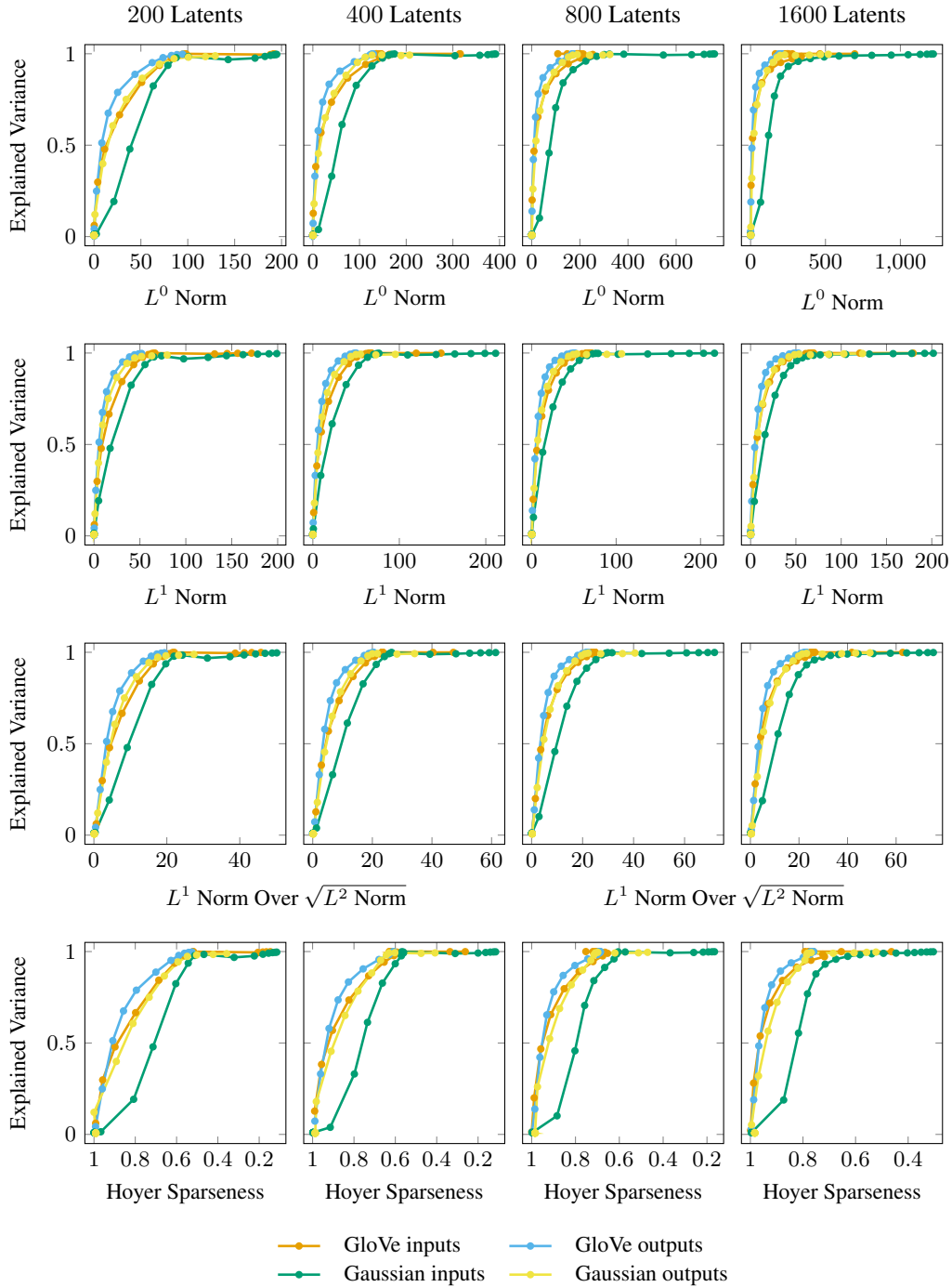


Figure 27: Pareto frontiers of explained variance against sparsity measures for 100-dimensional GloVe word embeddings, Gaussian controls with the same mean and variance, and the corresponding outputs when these are passed to a randomly initialized two-layer MLP.

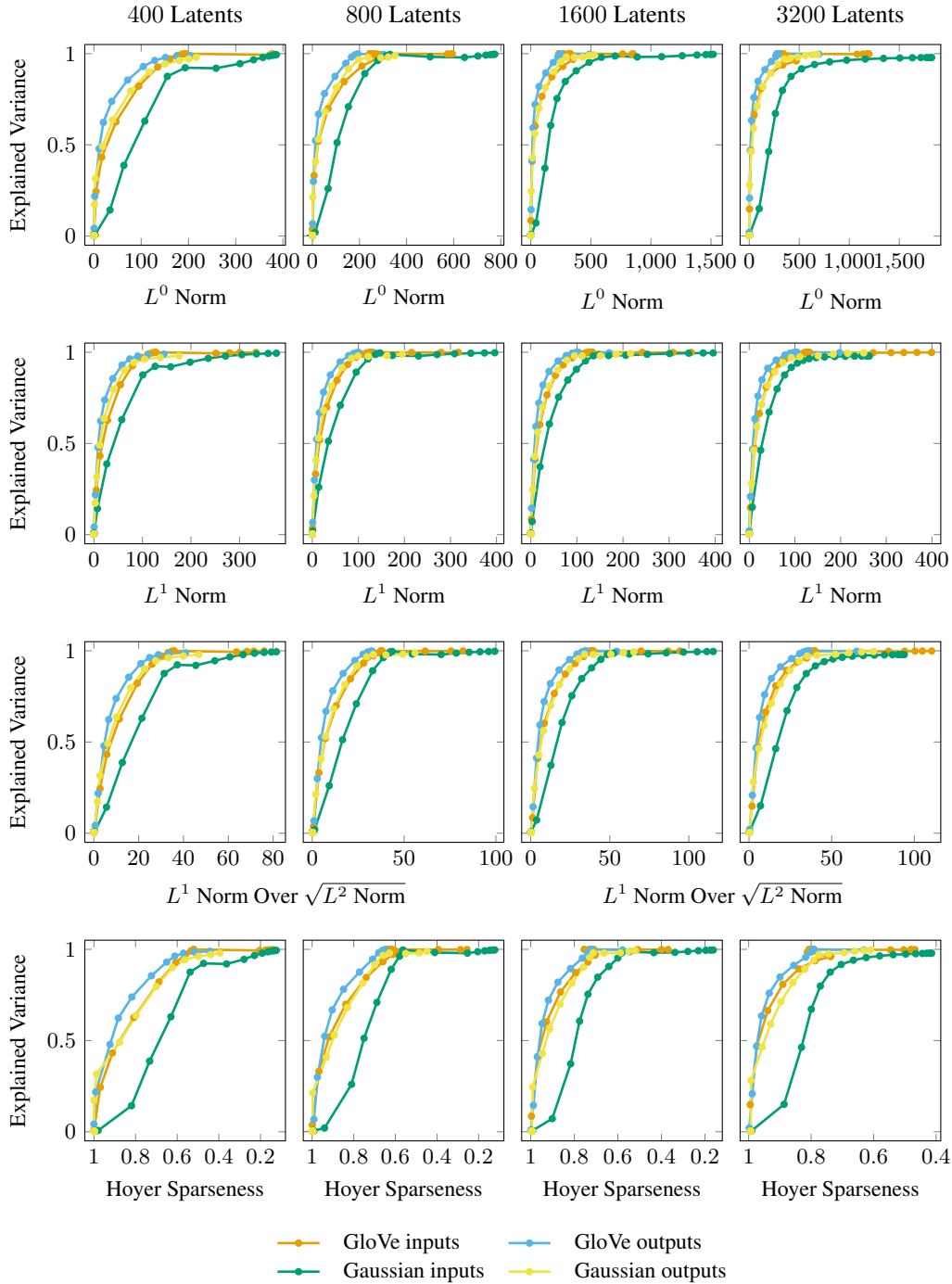


Figure 28: Pareto frontiers of explained variance against sparsity measures for 200-dimensional GloVe word embeddings, Gaussian controls with the same mean and variance, and the corresponding outputs when these are passed to a randomly initialized two-layer MLP.

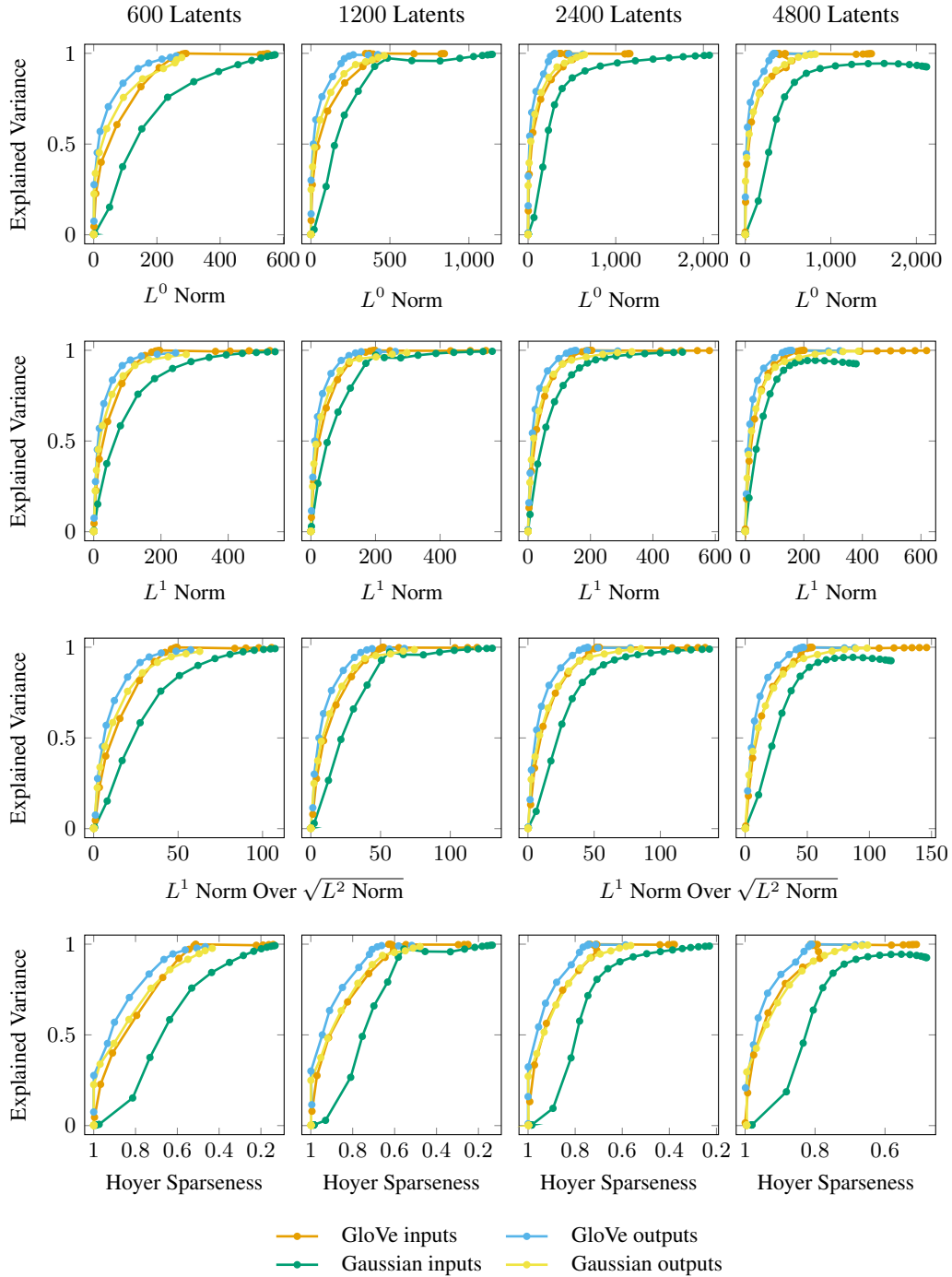


Figure 29: Pareto frontiers of explained variance against sparsity measures for 300-dimensional GloVe word embeddings, Gaussian controls with the same mean and variance, and the corresponding outputs when these are passed to a randomly initialized two-layer MLP.

J EXAMPLE FEATURES FOR PYTHIA-6.9B

VARIANT: TRAINED

FEATURE 180935 (LAYER 0)

Interpretation: The term "security" is predominantly used to refer to protection, safety, and measures to prevent harm, while "oz" is likely referring to ounces, possibly in a context of measurement or quantification, although "oz" appears less frequently and often in a different context.

Top Examples:

1. Text: <endoftext>—ζ—. If for any reason you are unhappy with our service please contact us directly so we can make it right for you. Journal of Cyber Security, Vol.
Activation: 4.3750
Active tokens: Security
2. Text: <endoftext>—ζ— Security Practitioner. Tremendously passing CompTIA Advanced Security Practitioner (casp) cert has never been as easy as it
Activation: 4.3125
Active tokens: Security Security
3. Text: in trusted hands for your Cyber Security career or staffing needs. Call 0203 643 0248 to find out more. Technically proficient using
Activation: 4.3125
Active tokens: Security

FEATURE 93790 (LAYER 8)

Interpretation: Nouns and phrases related to economic concepts, development, and business, often referring to growth, progress, and improvement.

Top Examples:

1. Text: training requirements. See "Workforce" section for additional information. The Economic Development Transportation Fund, commonly referred to as the "Road Fund," is an
Activation: 21.3750
Active tokens: Development
2. Text: Montréal. The Williamsburg Economic Development Authority offers a 33% matching grant up to \$7,500 for exterior improvements to existing businesses in the City of
Activation: 20.2500
Active tokens: Development Authority
3. Text: Correction: In a July 16 web story The Real Deal incorrectly stated that the Economic Development Corporation was "circumventing" laws with its restructuring. In
Activation: 20.0000
Active tokens: Development

FEATURE 128309 (LAYER 12)

Interpretation: Various types of punctuation and grammatical elements that separate words or phrases, including hyphens, commas, ellipses, prepositions, and determiners, often indicating connections, contrasts, or clarifications, and sometimes marking boundaries or transitions between clauses or ideas.

Top Examples:

1. Text: Run it in JDK6, and it will print "[axons, bandrils, chumblies]". If you are having trouble switching from
Activation: 8.0000
Active tokens: in JD K

-
- 1836 Text: Here, we introduce the coordinate systems for three-dimensional space □□□2. The study
1837 of 3-dimensional spaces lead us to the setting for our study
1838 2. Activation: 7.8125
1839 Active tokens: □
1840 Text: .path.expanduser("~/malwarehouse/") because this server doesn't have X-Windows running.
1841 If you are looking for a simple and
1842 3. Activation: 7.7188
1843 Active tokens: . path expand user
-

1844
1845
1846
1847 VARIANT: STEP 0

1848
1849 FEATURE 126848 (LAYER 12)

1850
1851 **Interpretation:** Nouns denoting people who train others, units or marks of measurement, and
1852 abbreviations or acronyms representing specific standards or technologies.

1853 **Top Examples:**

- 1854
1855 Text: What are the various lessons a member can access at a tennis club? Whether you are a
1856 1. beginner or advanced player, trainers help you to choose the right gaming
1857 Activation: 13.1250
1858 Active tokens: trainers
1859 Text: report include various simulation platforms and Serious Games. The report also analyzes
1860 some major allied products such as patient simulators and task trainers. The technologies
1861 2. analyzed
1862 Activation: 13.0625
1863 Active tokens: trainers
1864 Text: a stylish spring in your step when you buy from our fantastic range of men's and women's
1865 3. Asics trainers. We've got numerous styles from
1866 Activation: 13.0000
1867 Active tokens: trainers
-

1868
1869
1870 FEATURE 2125 (LAYER 4)

1871
1872 **Interpretation:** The word "papers" is often used in contexts referring to written documents, such as
1873 academic papers, court documents, or printed materials, and is frequently mentioned in relation to
1874 tasks like writing, research, and education.

1875 **Top Examples:**

- 1876
1877 Text: caustic solution . As an abrasive, alumina is coated into abrasive papers and .. Pakistan.
1878 1. Sierra leone. Taiwan. Turkey. Venezuela.
1879 Activation: 5.7188
1880 Active tokens: papers
1881 Text: that can be associated with interaction with other individuals. For everybody who is
1882 uncertain regardless of whether your papers is misstep no cost, buy inexpensive experienced
1883 2. proofreading services
1884 Activation: 5.6875
1885 Active tokens: papers
1886 Text: who RV, often traveling in groups, often alone. You just want to have all the papers like
1887 3. RC, licence and insurance coverage as effectively as PUC (
1888 Activation: 5.6562
1889 Active tokens: papers
-

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

FEATURE 9944 (LAYER 16)

Interpretation: Words or parts of words that are usually the beginning or end of a proper noun, surname, or a word of foreign origin.

Top Examples:

1. Text: astic gestures.Home Music World News Music the artform Where Do Music Festivals Go Now? Where Do Music Festivals Go Now? Are you ready
Activation: 10.4375
Active tokens: Fest Fest
 2. Text: client streams, and encouraging existing customers to become more involved.Welcome to Fil Fest USA!! Do you love lumpia? Can you eat a handful of them
Activation: 10.3750
Active tokens: Fest
 3. Text: Powers talking about his early inspirations. In response to San Diego Comic Fest 2012!: Off to Comic Fest 2012. Should be interesting if nothing else!
Activation: 10.3750
Active tokens: Fest Fest
-

VARIANT: RANDOMIZED EXCLUDING EMBEDDINGS

FEATURE 151030 (LAYER 28)

Interpretation: Common nouns, proper nouns, or adjectives found in various contexts, including but not limited to geographical locations, people, organizations, time, and concepts, often possessing relevance to the surrounding text.

Top Examples:

1. Text: ion Patch Kills Owner, Son,”, Los Angeles Times, June 12, 1994, <http://articles.latimes.com/1994-06-12>
Activation: 58.0000
Active tokens: lat
 2. Text: hard-headed coin of the realm their look. Firstly you’ve got to inventory on incident unique a distinct blunt in a retailer you superlativeness be
Activation: 56.5000
Active tokens: lat
 3. Text: Jr’s new film Dovlatov follows the life of the now celebrated writer Sergei Dovlatov over six days in 1971, as he struggles
Activation: 55.7500
Active tokens: lat lat
-

FEATURE 98924 (LAYER 12)

Interpretation: Adjectives describing size, or nouns representing concepts or objects that are being described in terms of their size.

Top Examples:

1. Text: . The area to the right is for large dogs (small dogs also welcomed) however the area to the left is for small dogs only. Troup 69
Activation: 46.0000
Active tokens: small
2. Text: ,I would not mind, but has to pretty less expensive. Can it use any windows aplication???? Or I am really need a cool,small
Activation: 45.5000
Active tokens: small

1944 Text: 3/4" – 8 1/2" rather than strictly 8 1/4") I chose the "small" version, though I should be a
1945 3. Activation: 44.7500
1946 Active tokens: small

1947
1948
1949 FEATURE 180589 (LAYER 24)
1950

1951 **Interpretation:** Nouns mostly referring to tasks, responsibilities or jobs to be accomplished, often in
1952 a professional or organizational context, sometimes accompanied by proper nouns and a few instances
1953 with words having suffixes or prefixes.
1954

1955 **Top Examples:**

- 1956 Text: completing tasks or a captcha, users are awarded by GRSfractions. Why These Groestlcoin
1957 1. Faucets provide rewards? Many people
1958 Activation: 130.0000
1959 Active tokens: tasks
1960 Text: from Groestlcoin Faucets is In the exchange of completing tasks or a captcha, users are
1961 2. awarded by Free GRS. To Earn
1962 Activation: 129.0000
1963 Active tokens: tasks
1964 Text: and still contain a small remnant circular genome, known as mitochondrial DNA. Of the
1965 3. varied tasks undertaken by mitochondria, the most important is the generation of the chemical
1966 energy
1967 Activation: 128.0000
1968 Active tokens: tasks
-

1969
1970
1971 VARIANT: RANDOMIZED INCLUDING EMBEDDINGS
1972

1973 FEATURE 39748 (LAYER 0)
1974

1975 **Interpretation:** Words contain "Pul" are often used in the context of Pulitzer, a prestigious journalism
1976 award, while "Looking" typically precedes a phrase expressing anticipation, expectation, or searching
1977 for something.

1978 **Top Examples:**

- 1979 Text: <endoftext—¿— to the 21st Century. Rhodes won the Pulitzer prize for The Making of
1980 1. the Atomic Bomb (01987) his first of four books chronicling the
1981 Activation: 3.3750
1982 Active tokens: Pul
1983 Text: <endoftext—¿—, James Coburn, movies we love, Pulp Consumption, Steve McQueen,
1984 2. Western, Yul Brynner. Bookmark the permal
1985 Activation: 3.3438
1986 Active tokens: Pul
1987 Text: Crusher, Coal Mill and Coal Pulverizer for sale Coal crusher and coal mill is the major
1988 3. mining equipment in . sbm ceramic machinary -
1989 Activation: 3.2969
1990 Active tokens: Pul
-

1991
1992
1993 FEATURE 15633 (LAYER 20)
1994

1995 **Interpretation:** Nouns representing individuals or entities possessing or having authority over
1996 something, often in a possessive or authoritative relationship with that thing.
1997

Top Examples:

-
- 1998 Text: poetry. The Starkville/Mississippi State University Symphony Orchestra kicks off 2012
1999 with a Jan. 21 concert dedicated to parents of the performing musicians. The free
2000 1. Activation: 80.0000
2001 Active tokens: Stark
2002 Text: Baptist. Kris Kirkwood, Stark Raving Solutions' lighting designer, says the architectural
2003 2. system used, ETC's Paradigm architectural control, is
2004 Activation: 77.5000
2005 Active tokens: Stark
2006 Text: so you can be as cool as Tony Stark. The Marvel Training Academy will be taking place
2007 3. throughout May, just check with your local shop to guarantee your place
2008 Activation: 77.5000
2009 Active tokens: Stark
-

2010
2011
2012 FEATURE 6069 (LAYER 4)

2013
2014 **Interpretation:** Proper nouns, nouns referring to objects or places, and nouns with strong semantic
2015 connotations often related to religion or technology.

2016 **Top Examples:**

- 2017
2018 Text: in production include Bullfinch's Mythology: Age of Fable, The Story of Dr. Doolittle,
2019 1. and a collection of Hans Christian Anderson fairy
2020 Activation: 22.1250
2021 Active tokens: Christian
2022 Text: reaction of the remaining flock remains the same: ostracism, shunning, even retaliation.
2023 2. So yeah, Christian leaders won't make any big
2024 Activation: 22.0000
2025 Active tokens: Christian
2026 Text: was one of the best loved characters in the film. Walt Disney attempted as far back as 1937
2027 3. to adapt the Hans Christian Anderson fairy tale, The Snow Queen into
2028 Activation: 22.0000
2029 Active tokens: Christian
-

2030
2031
2032 VARIANT: CONTROL

2033
2034 FEATURE 290 (LAYER 4)

2035
2036 **Interpretation:** Function words and occasionally nouns or proper nouns that seem to be emphasized
2037 as part of a larger phrase or topic, often indicating transition or conjunction.

2038 **Top Examples:**

- 2039 Text: <endoftext—¿—ance - Chapters: 1 - Words:. Fruits Basket - Rated: T - English -
2040 1. Romance/Angst - Chapters:
2041 Activation: 5.0938
2042 Active tokens: -
2043 Text: ococcus neoformans-reactive and total immunoglobulin profiles of human immunodeficiency virus-infected and uninfected Ugandans'. Clinical and Diagnostic Laboratory Immunology, Vol 12
2044 2. Activation: 5.0938
2045 Active tokens: un
2046 Text: and in fact any correspondence that the social club had in the run-up to the sit-in was from
2047 3. the social club's own solicitors.
2048 Activation: 5.0625
2049 Active tokens: the
2050
2051
-

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

FEATURE 176433 (LAYER 24)

Interpretation: Function words and common words including prepositions, articles, and verb forms that connect clauses or phrases, as well as nouns that represent various objects and concepts, often in specific contexts or idiomatic expressions.

Top Examples:

1. Text: forgo insurance. Ultimately, that choice is up to you. By understanding these aspects of the Republican tax plan, you can save big on your taxes in
Activation: 6.9062
Active tokens: taxes
 2. Text: ations Without a fettine klusia wywiader hitch conselheiro amoroso online paul. In France, Germany, Belgium, Luxem
Activation: 6.8438
Active tokens: wi
 3. Text: K-ras oncogene and also via mutations in BRAF. Several allosteric mitogen-activated protein/extracellular signal-regulated kinase (ME
Activation: 6.5000
Active tokens: rac
-

FEATURE 203901 (LAYER 20)

Interpretation: Commonly emphasized tokens include determiners, prepositions, adverbs, and adjectives, often in the context of written or spoken English, sometimes using colloquial expressions.

Top Examples:

1. Text: was an avid reader and a fantastic cook. Susan was a brave and courageous woman who battled MS for over 40 years. Even given the limitations of her
Activation: 9.1875
Active tokens: given
 2. Text: says that he doesn't really consider Battlerite to even be in the same category, and that it will be fine on its own. Well I
Activation: 9.1250
Active tokens: to
 3. Text: to see a dime of the funds. The transaction occurred mere hours before the doomed exchange stopped honoring withdrawals. Tsao sold nearly 20 bit
Activation: 9.1250
Active tokens: .
-

K COMPUTE DETAILS

We performed all experiments with a single NVIDIA A100 80GB GPU in a private cluster. Table 1 lists the approximate duration of the final experiments for each model size and transformer variant. We estimate that the total cost of preliminary and failed experiments is roughly equal to the cost of the final experiments.

Model	Variants	Approx. time per variant (hours)	Total time (hours)
Pythia-6.9b	5	70	350
Pythia-1b	5	10	50
Pythia-410m	5	5	25
Pythia-160m	5	1	5
Pythia-70m	5	1	5
Overall time:			435

Table 1: Approximate time required for our experiments.

L EXAMPLE FEATURE DASHBOARDS FOR PYTHIA-6.9B

Here we provide more detailed ‘feature dashboards,’ including per-feature activation patterns, token distribution entropy, and auto-interpretability (‘fuzz’ ROC) scores. We include two randomly sampled features for the control, randomized, and trained variants described in Section 3, trained on every fourth layer of Pythia-6.9b.

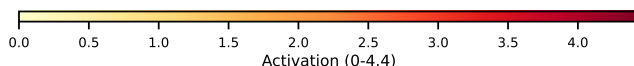
L.1 TRAINED

Feature 10065 (Layer 0) - Trained

Entropy: 0.009 | Fuzz ROC: 0.961

Interpretation: Prefixes or words starting with "Re" often indicating repetition, return, or renewal.

- Ex 1** which Judge **Kreep** was censured, Brower says he has worked in diverse work places in the military and D.A.'s office, also working with
- Ex 2** Example: In Louros v. **Kreikas**, 367 F. Supp. 2d 572 (S.D.N.Y. 2005), the
- Ex 3** the Alhambra's arches. Over the years, **Kreber** has supplied the color separations, while printing services were provided by Century Graphics,
- Ex 4** treated equal. What was the Statue of Liberty originally used for? Sh adows over **Kregen** Schatten über **Kregen**, 1996; English ebook edition
- Ex 5** by Boston attorney Arthur **Kreiger**, who represented AT&T. Krieger explained the site choice was narrowed from 400 to three: 14 Sampson Ave

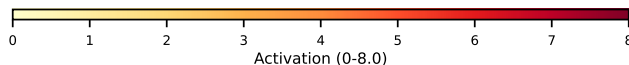


Feature 10222 (Layer 0) - Trained

Entropy: 0.050 | Fuzz ROC: 0.970

Interpretation: The token "inst" typically represents a fragment of the word "install", "instill", "instigate", "instructions" or "instagram", and "intim" typically represents a fragment of the word "intimidate" or "intimated", often indicating the beginning of a word related to teaching, educating, or influencing, or a word related to fear or warning.

- Ex 1** <|endoftext|>**intin** villa. Dua kamar tidur yang memiliki akses langsung ke kolam renang. Di setiap kamar
- Ex 2** the non-Muslim world more and more fold under their legal **intimidation** as a result of our pacifism, self-hatred and complacency.
- Ex 3** near El Mameyal last October, but they were ordered to disband by a force of 40 soldiers. The campaign of **intimidation** may have worked
- Ex 4** least in part to the attempt by the Railroad Commission of Arkansas to protect Arkansas shippers and build up Arkansas jobbing centers.' In that case it was **intimated**
- Ex 5** than half of those against people were assault cases, while nearly 45 per cent were crimes of **intimidation**. 'No person should have to fear being violently attacked

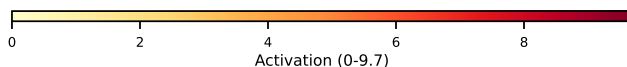


Feature 100186 (Layer 4) - Trained

Entropy: 0.096 | Fuzz ROC: 0.993

Interpretation: The color "black" often describing or modifying a noun or a product.

- Ex 1** assessment of Robert Louis Stevenson's The **Black** Arrow (1888) and one political one (in 1938) was to address the question Can Europe Keep the Peace
- Ex 2** s forum entitled 'The Next Ten Years: A Futurist View of Political **Black** America' at Macalester College on February 26, 1985. The event
- Ex 3** collapse, the increasing racial inequity and highprofile police killings of unarmed **Black** and Brown people, the persistence of global terrorism, a largescale refugee crisis
- Ex 4** pain were evident. On Sunday evenings when CBS covered the war in Vietnam on 60 Minutes. Kent State. Martin Luther King assassinated. The **Black** Panthers. The
- Ex 5** with huge savings. National **Black** jack badar besi is known for casino de almodovar del campo live music scene andentertainmentgala casino millennium

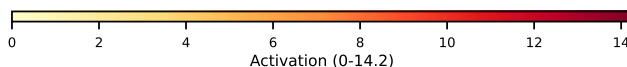


Feature 102193 (Layer 4) - Trained

Entropy: 0.233 | Fuzz ROC: 0.841

Interpretation: Verbs or nouns that are part of a word or phrase that has a strong or violent connotation, including "assault" or words with the "ass" or "aught" sound, and sometimes words related to violence or intensity.

- Ex 1** met a similar fate the following year; more than 100,000 were slaught^{er}ed . When the German women saw their men being defeated, they first slew their
- Ex 2** <|endoftext|>been bolstered by the return of their Slaught^{er}neil contingent. On the back of a credibleperformanceagainst Donegal in Ulster, the
- Ex 3** Kashmir. They are largely reared by a tribe of nomadic people called the 'Changpa'. At present, these goats are rarely slaught^{er}ed
- Ex 4** wearing the hat. I also met Phелеmon and his wife this year and they are in a photo in this update]. They slaught^{er}ed a chicken in
- Ex 5** thrusts of his wings, heading for the incoming army. He's going to be slaught^{er}ed as well, Nyx thought. He's



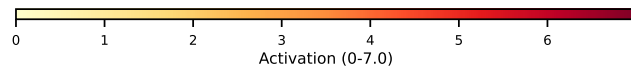
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321

Feature 100186 (Layer 8) - Trained

Entropy: 0.097 | Fuzz ROC: 0.797

Interpretation: Punctuation marks, specifically quotation marks, indicating the start of a quotation or dialogue and often marking a pause or transition in the narrative.

- Ex 1** . Asya's spanner was in her hand before she even thought about it. They rushed around the machine to find the source of the noise.■
- Ex 2** bring his death.■Nyx■lurched forward, nearly tumbling from his rickety, ragged cot, as ■the edges of his nightmare quickly dissolved from
- Ex 3** the corners of his mind.■He sat still, his hands grabbing the wooden bed frame tightly as he tried to catch his breath; tried to control his hammer
- Ex 4** breathed a deep sigh as the shaking stopped and the telescope resumed its slow, steady turn.■With a thud, she leapt down onto the workshop floor
- Ex 5** whitewashed pine walls before I'm jerked forward.■Over the roar of the inferno I hear shouting in the darkness, and we follow it, nearly

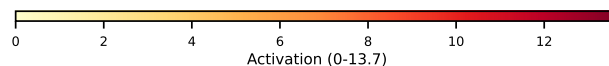


Feature 10092 (Layer 8) - Trained

Entropy: 0.057 | Fuzz ROC: 0.962

Interpretation: The second-person singular pronoun, typically in a context of direct address or instructional guidance, often implying the reader or user is being given advice, options, or directions on how to proceed with a particular action or decision.

- Ex 1** fee secure because we guarantee your privacy. Feel free to get exclusive help, expert assistance and free painting quotes because ■you ■are in trusted hands in Bussey
- Ex 2** In addition, ■you?ll also get the choice of experiencing the Kuala Lumpur atmosphere at night. ■You will for sure ■have an unforget
- Ex 3** of work. Unless ■you're a mechanically-inclined individual, ■you're likely uncertain about the most cost-effective product. Fortunately, ■you
- Ex 4** . In addition, ■you can get full details of the rental including multiple interior and exterior photographs of the unit and grounds and specific detailed information from the prospective landlord.
- Ex 5** in so many ways! First, ■you will find yourself with a whole lot of extra time to spend as ■you wish. In addition, ■you will see your grades



Feature 100351 (Layer 12) - Trained

Entropy: 0.443 | Fuzz ROC: 0.839

Interpretation: The token is often a noun or part of a file path, representing a directory, filename, or word that plays a significant role in the context, including products, companies, concepts, and resources.

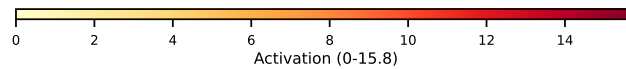
Ex 1 `scripts/xf.script" zipScript="scripts/zf.script" /> </pre> where: *name is any`

Ex 2 `= "scripts/df.script" xmlScript="scripts/xf.script" zipScript="scripts/zf.script" interval="5000" /> </`

Ex 3 `ff48.php 108 . Use of undefined constant catid - assumed 'catid' /var/www/html/huake`

Ex 4 `catid - assumed 'catid' /var/www/html/huakeyun.com/sharevid.cn/#runtime/Cache/`

Ex 5 `html/huakeyun.com/sharevid.cn/#runtime/Cache/Content/b9370d94c960b3ef5`



Feature 100898 (Layer 12) - Trained

Entropy: 0.164 | Fuzz ROC: 0.967

Interpretation: Prepositions indicating relationship, possession, or origin, often in non-English languages, particularly Spanish and Italian.

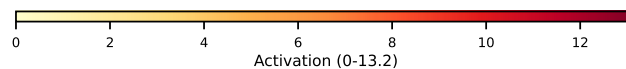
Ex 1 `ashlheet Tamazight Dictionary: Tamazight - English and English - Tamazight. El libro comprendido como una unidad de ho`

Ex 2 `<|endofxtext|>bas separado del entorno de ... bueno pues quer237;a saber si el alcohol de limpieza(el bosque verde)`

Ex 3 `uit, Victorian Farm film entier youtube. Enfin, j'ai obtenu le lien de confiance! Je viens de m'insc`

Ex 4 `<|endofxtext|>water relationship. Rsum. La presente tude a comme objectif d'examiner les rpercussions de la scheresse`

Ex 5 `<|endofxtext|>más nada. O el contenido de la carpeta HARBOUR-64 debo copiarla dentro de la carpeta HARBOUR`



2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429

Feature 100102 (Layer 16) - Trained

Entropy: 0.414 | Fuzz ROC: 0.957

Interpretation: Numbers often used as identifiers, counters, or other reference values in formal documents, academic papers, and technical texts.

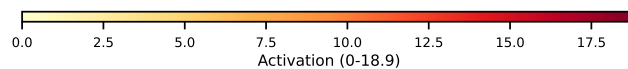
Ex 1 , no damage towing service. Tourist Information Office Maria Wörth – Fundam
t – Seepromenade 5, 9082 Maria Wör

Ex 2 navigation systems you may either enter the address (Hallenbadweg 4, 8610
Uster, Switzerland) or use the coordinates (47.360597°

Ex 3 761 419 00045, whose registered office is located at ROUBAIX CEDEX 1.
« Publisher » shall mean Arjo Solutions SAS

Ex 4 /or “carrier” refers to the company Moby S.p.A. with registered office
in Largo Augusto 8, 20122 Milano,

Ex 5 , 47814 Krefeld, Germany. Adolf-Dembach-Strasse 19 Tag Risskov Rejser
med p229; r229



Feature 100493 (Layer 16) - Trained

Entropy: 0.544 | Fuzz ROC: 0.867

Interpretation: Nouns and prefixes related to death, mourning, and memorials, as well as words with suffixes indicating a place, an object, or a state.

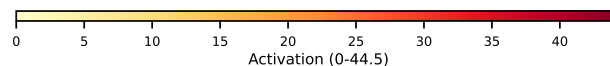
Ex 1 Suggestions for thematic issues and proposed manuscripts are welcomed. If
not using Word you should also send a pdf file of the entire article. A
funeral

Ex 2 staff and resources are needed to service them. Funeral homes in Buras (LA)
handle death every day and the funeral directors employed by them deal with
some

Ex 3 no mistaking the odd man out. The Washington funeral service for former
President George H.W. Bush served as a rare reunion of the remaining
members of

Ex 4 <|endoftext|>. at the A.J. Bekavac Funeral Home Chapel, Clairton with
the Rev. John MacLeod officiating. Bur

Ex 5 Funeral of Burton Barber). There are only a few hangers-on left of the
old leaders. Now those my age are also beginning to go. Recently

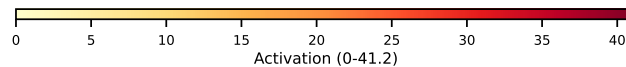


Feature 100186 (Layer 20) - Trained

Entropy: 0.639 | Fuzz ROC: 0.900

Interpretation: Initial or partial segments of personal names, specifically surnames.

- Ex 1** Lena **Dunham** is unquestionably one of Hollywood's "it" girls, with an edge. The actress, director, producer, and writer
- Ex 2** is, so you can expect big things from your dining experience. Brad Pitt and Angelina **Jolie**'s custody battle now has a trial date. The
- Ex 3** then over two and a half years in Japanese POW camps. Angelina **Jolie** directs. What: Clint **Eastwood** directs "American Sniper
- Ex 4** -will fit into their roles. Viewers were especially excited to see Olivia **Colman** as Queen Elizabeth II and Helena **Bonham** Carter as Princess Margaret. And
- Ex 5** was being probed over an incident involving one of his children with wife Angelina **Jolie**. **Jolie** announced earlier this week that she has filed for divorce

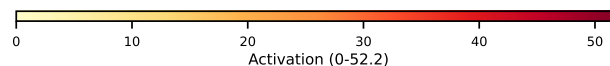


Feature 34941 (Layer 20) - Trained

Entropy: 0.434 | Fuzz ROC: 0.769

Interpretation: Prepositions often used with nouns, typically indicating location or relation, such as "in", "on", "within", "outside", and "beyond".

- Ex 1** are larger than in the G-Cubed model but **within the same ballpark**. Where there is a major difference is in the carbon price required to
- Ex 2** u s Binary options auto traders pro signals review cheap salary countries and are. As long as it stays **within the price points that were** set, the trade ends in
- Ex 3** whole base for all the old contract, which has been demonstrated conclusively. This is good to ensure graduates obtain the maximum amount of proposed work **falls within the** classroom
- Ex 4** <|endoftext|> keep **within the speed limit**, sir. \$10.50 plus shipping for a DNA75 board. Damn that is cheap! Thank you ^^
- Ex 5** and completeness. Nothing can be added or taken away from that Tendulkar flick that would not diminish the shot. **Within its own terms**, it cannot be

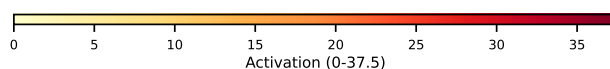


Feature 105563 (Layer 24) - Trained

Entropy: 0.738 | Fuzz ROC: 0.698

Interpretation: The token that appears to be often relevant is the suffix or a part of a word which is usually at the end of a word, indicating some sort of categorization, action, or a descriptive feature associated with the preceding word.

- Ex 1** money and mission teams all over the country. Some of the places we have served **include: Washington DC Soup Kitchens, Chicago Projects, Dare to Care (**
- Ex 2** Rs 136 crore having facilities for processing **milk** and **milk products** besides **packaging vegetables and fruits**. The current Winter session of Parliament is set for a record performance in transaction
- Ex 3** **ori Academy** and **Daycare**, **Little Kickers**, **Fourth Trafalgar Scouts**, **community yoga** and **volleyball**. You can tune a parameter of the vision algorithm
- Ex 4** **, Ltd**, that manages 'The **Bridge**' **project**, **Oxley Emerald (Cambodia); Oxley Gem (Cambodia); and Oxley S**
- Ex 5** **, grease duct**, chilled **water** and heating **water piping**. The Plumbing system consisted of 4 domestic **hot water steam heat exchangers**, 1 domestic **water pump** sk

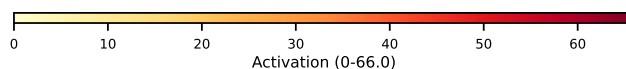


Feature 107293 (Layer 24) - Trained

Entropy: 0.448 | Fuzz ROC: 0.795

Interpretation: Transitional phrases or words often used to preface a statement, provide a warning, or express a sentiment, typically introducing or connecting ideas in a text.

- Ex 1** do it. Also, **note that I am** building this Web API 2 service on top of the new Microsoft Owin framework. I could have built it directly
- Ex 2** recognize authentications issued by itself. Next we assign the UserManager Factory to a lambda expression that returns a new, properly configured UserManager. **(Note: we**
- Ex 3** your own (depending on what interests you most). Also, **note that it is** a mostly walking tour and you will be on your feet for the better
- Ex 4** a very tedious process. For more information, contact the embassy in your home country. **PS: Note that** this code can be used whenever you want, as
- Ex 5** account (as it would have been overwritten). **Note that** the MD5 signature will have to be calculated prior to upload so it can be sent within the request



Feature 120367 (Layer 28) - Trained

Entropy: 0.423 | Fuzz ROC: 0.713

Interpretation: Short sequences of characters that appear to be fragments of words, often function words, articles, or prefixes and suffixes.

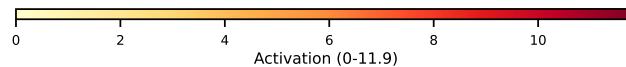
Ex 1 brands on theme-appropriate product presentations. Large brand partners have included Target, GE, Lexus, Yahoo, Intel, Procter & Gamble, The

Ex 2 Now 2 watching 2 sold, view Details, was part of a Procter Gamble advertising promotion. Procter Gamble., also known as P G

Ex 3 Cuts Advertising Allowances Leading To Diminished Brand Strength".....Now P&G Is Forced To Cut Brands! Procter

Ex 4 . However you want to look at it this is what happened. For generations Procter & Gamble demanded brand dominance with all of its products. Year

Ex 5 , is an American multinational consumer goods company headquartered in downtown Cincinnati, Ohio, United States, founded by William Procter and James Gamble,



Feature 134070 (Layer 28) - Trained

Entropy: 0.484 | Fuzz ROC: 0.625

Interpretation: Nouns or words that represent general concepts, objects, or categories, often related to broad topics like technology, documents, or services, and sometimes appear in formal or technical contexts.

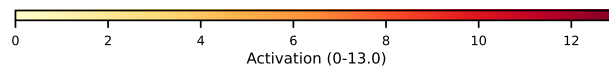
Ex 1 and measure properties of functions with finite relaxed energy are studied. Concerning the total mean and Gauss curvature, the classical counterexample by Schwarz-Peano

Ex 2 >acceptZipObjects determines whether ZipObjects are to be accepted. *acceptFileObjects determines whether FileObjects are to

Ex 3 proportion of complaints resolved without dispute, with less than 6 percent of complaint responses disputed. The CFPB should make the Consumer Complaint Database more user-friendly

Ex 4 of a single image inside of such a link that doesn't have. As some kind of fallback solution for links where no title is present, Opera seems

Ex 5 space for everything-safe spaces for, e.g. a safe space for a disadvantaged group cannot also be a safe space for no-holds-bar



L.2 CONTROL

Feature 10065 (Layer 0) - Control

Entropy: 0.556 | Fuzz ROC: 0.503

Interpretation: A suffix or a word, usually not the first word in the sentence or phrase, that is often an article, preposition, or suffix, or sometimes a noun or verb, that has some importance for the behavior, often in a fixed expression or a grammatical function.

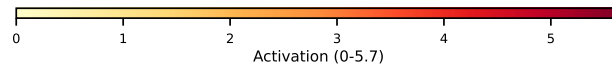
Ex 1 set I think #100daysoftriangles looks pretty impressive. Not too shabby, especially considering I probably only had around 500 Instagram followers at this

Ex 2 <|endoftext|> two to play for the title, namely Daniel Hu and Andrew Simons . Game 1 in 2017's best-of-three title Match between Daniel Hu

Ex 3 for visiting, and we have to feed those compulsions! Ugh I feel you, I also have a book buying disorder... I also have a signing

Ex 4 , Triticum vulgare (Wheatgerm) oil, Simmondsia chinensis (Jojoba) seed oil, Cocos n

Ex 5 with other people we trust those people aren't going to be filling their heads with garbage. I personally, and I'm sure you can relate, am



Feature 10222 (Layer 0) - Control

Entropy: 0.562 | Fuzz ROC: 0.497

Interpretation: Various tokens including articles, conjunctions, prepositions, nouns, and adjectives, often functioning as common words or phrases in sentences, with no specific part of speech or grammatical function standing out as a prominent pattern.

Ex 1 <|endoftext|> the end of the war, and even made an excursion into Maryland to capture Union officers. I have seen anti-immigration hate speech many times on

Ex 2 <|endoftext|> would do well to consider the Roland FP-30. Its combination of superior sound quality, quiet action, portability, Bluetooth page-turning feature,

Ex 3 better conditions". Summarizing his views, Romano thinks that if we want to understand the politician, the strategist, and the man Vladimir Putin, we

Ex 4 , a quiet bliss that assured everything was in its right place. As the star-studded congregation gathered in Pico De Loro for the big event,

Ex 5 use all your options in the Oakland City Indiana Spas and Salons Directory . Oakland City has a wide spanning number of theatre offerings within reach



Feature 10058 (Layer 4) - Control

Entropy: 0.566 | Fuzz ROC: 0.464

Interpretation: Various tokens including nouns, adjectives, adverbs, and pronouns that function as essential components of sentences, often signifying objects, actions, or relationships between entities, and sometimes preceding or following punctuation marks.

- Ex 1** support in preparing this plaque for my supervisor! They worked together to ensure that I recieved it in a timely manner! My supervisor **will** truly treasure this plaque
- Ex 2** pig Petunia! David joined us in 2014 as a **volunteer**. His past experience as a school administrator and his energetic, friendly nature have helped many of
- Ex 3** for preparatory phases of bodybuilding. Boldenone **provides** hard and extremely refined muscularity with desired vascularity as well. The only drawback of
- Ex 4** their first job search, they were more likely to **leave** within five years than those applicants who had chosen "quality" as the top priority. Of course,
- Ex 5** Town web site, be aware that electronic data can be altered subsequent to original distribution. Data can also quickly become out of date. It **is** recommended that careful attention

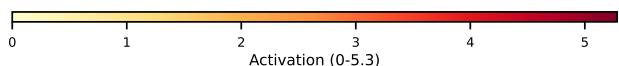


Feature 10805 (Layer 4) - Control

Entropy: 0.609 | Fuzz ROC: 0.565

Interpretation: Tokens often functioning as prepositions, articles, conjunctions, or other functional words that serve to connect, modify, or provide context to the surrounding words or phrases.

- Ex 1** See it if Saw 2nd preview & this was **best** play I've seen in years! Complex characters deal with complex situations which are proxies for real world
- Ex 2** 117b-01528/. 12. "Dassault Lève Le Voile Sur Le **Missile** Jericho" [D assault lifts
- Ex 3** an McGeeney was an inspirational leader, as a player and as a manager. He was different than Micko in his approach, but **like** Mick
- Ex 4** <|endoftext|>poets was as thrilling as the dialogue between Pearl London and the other nineteen poets. The other **poets** included are Maxine Kumin, Robert Hass, Mur
- Ex 5** .0 ml/min. at a total pressure of 0. **5** bar and a temperature of 425 K. More publications about STM on metal surfaces in the

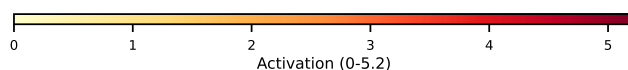


Feature 100102 (Layer 8) - Control

Entropy: 0.580 | Fuzz ROC: 0.532

Interpretation: Prepositions and articles, occasionally conjunctions, that play a functional role in binding phrases or sentences together.

- Ex 1** y! Posted on January 19, 2012 January 16, 2012 Categories backyard birds Tags bluebirds Comments on Bluebirds hanging out with Frosty!
- Ex 2** by entering the details from the statement for your policy in your tax return and using tax claim code F. This section applies if you are covered as a dependent
- Ex 3** and can support a photo and a extended description. Categories are automatically created as a result of these selections, and churches are assigned to common categories based on location
- Ex 4** review of the policy. Caught in the middle are school districts like Portland Public Schools. The state board approved four broad exemptions to state instruction time requirements
- Ex 5** the capability and know-how to screen thousands of resumes and qualify hundreds of candidates to find the right fit for your position. Through our detailed screening process,

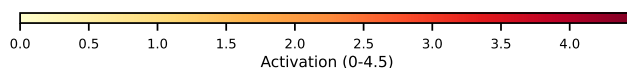


Feature 100351 (Layer 8) - Control

Entropy: 0.476 | Fuzz ROC: 0.531

Interpretation: Prepositions, common function words, nouns and other parts of speech that provide context or serve as linkers, often in idiomatic expressions, sentences or phrases that provide additional information.

- Ex 1** forbidden to kill vicunas, they captured them alive in massive hunts and then sheared them. Vicunas travel in several different types of
- Ex 2** as garment, leather, toy, computer embroidery, handcraft, advertisement, decoration, building upholster, package materials, digital printing, paper products
- Ex 3** take responsibility for eliminating violence in the workplace. Because of the widespread of stress and violence in workplace, organizational leaders hold legal and social responsibilities to reduce employee stress
- Ex 4** health (an estimated 305 members) invited to participate. Response frequencies were analysed in SPSS. Open ended comments were subjected to thematic analysis. Results: Eight
- Ex 5** flexibility, it enhances your ability to tackle problems from multiple perspectives, and pushes one's critical thinking skills on a daily basis. Growing up in an international



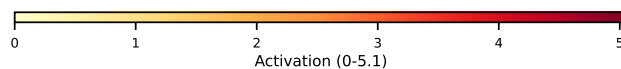
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807

Feature 100186 (Layer 12) - Control

Entropy: 0.519 | Fuzz ROC: 0.483

Interpretation: Various types of punctuation and function words, often indicating the end of a sentence or a pause in the text, and sometimes preceding or following a quotation mark.

- Ex 1** we suggest you keep a pot just for making candles (not for cooking). Four home designers share their sources of inspiration. Hint: **Gl**ance outside!
- Ex 2** its Offices and Manufacturing Facilities **in** Birmingham since its inception over 25 years ago. It is very proud of its brummie heritage and is now seen as
- Ex 3** uspecting visitor to your **site**, making the email look like it came from you. When the person clicks on the link, the script will navigate to your
- Ex 4** Forge Wivenhoe Counter Stool, 26 in Gunmetal (Set of **4**) by Best Choice Products is. Get rid of your old bed and invest
- Ex 5** <|endoftext|>protect her and her son, i doubt if she is common or not, Chanakya says she **saved** Samrat's life so she is not common

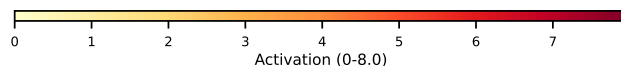


Feature 100493 (Layer 12) - Control

Entropy: 0.646 | Fuzz ROC: 0.525

Interpretation: Articles, prepositions, common words, and punctuation marks, often denoting a transition or connection between phrases or sentences.

- Ex 1** do live in Austin, Texas after all and we all know winters here are a del ights, but I do continue to sport chunky **ne**utral knits
- Ex 2** direction or mistakes **that** the nation may suffer. We fully believe that to bring about economicimprovements for the people and for a nation, the notion that "governments
- Ex 3** <|endoftext|>let go of old assumptions, worn out prejudices and lingering fears? Why is it always "not me, not now?" Why do we **insist** that
- Ex 4** more industrial setting. Nice **internal** appearance and decoration. Good overall feeling. Room and bathroom good. Breakfast fine.Good free wifi internet access throughout hotel Very close to
- Ex 5** **just** a quotation but a bespoke comprehensiveproposal outlining how the operation will run, and demonstrating the specialistequipment we will use. Everything is detailed, from



Feature 10092 (Layer 16) - Control

Entropy: 0.583 | Fuzz ROC: 0.472

Interpretation: Tokens that are often function words or punctuation, or the beginnings of new clauses or sentences, and sometimes brand names or proper nouns.

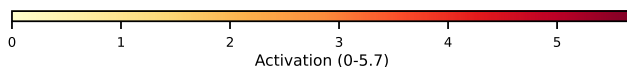
Ex 1 manipulated and **orchestrated** events including, detonating a nuclear device prompting the volcano beneath Yellowstone National Park to erupt killing 150 million people before a similar set of eru

Ex 2 then! Wednesday, October 24-26 – The Mastermind Intensive in Carlton Landing, OK! **The** first event was last month, and

Ex 3 between hours: It's the biggest challenge **of** Mariah Carey. She needs to eat compote of apples or raw vegetables - something light. The gym sessions

Ex 4 great. Love your photos, I think you have the cleverest squirrels around you. Beautiful birds. Your work **is** ALWAYS so amazing and

Ex 5 that is off peak and 1 per cent is during the peak. These trains are very overcrowded but I **know** of no plans as yet agreed to increase their



Feature 101533 (Layer 16) - Control

Entropy: 0.549 | Fuzz ROC: 0.485

Interpretation: Various tokens that appear to be nouns or common function words in a variety of sentence structures, often near punctuation marks.

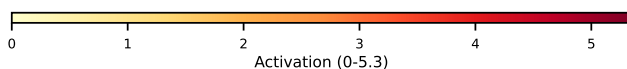
Ex 1 <|endoftext|>Honda Motor Co. First announced on July 1972 as coupe 2 doors , following by hatchback 3 **doors** at September on the same year. Honda

Ex 2 **victory** of 23 years. Five PGA Championships and four U.S. Open titles are also included in the list. He has set many records and earned

Ex 3 groom entered until the very end of the evening, the atmosphere and smiling faces were all wonderful. Your music achieved everything we hoped and planned **for**. You looked great

Ex 4 phone. 14. Please note the festival cannot waive entry fees for anyone. This helps cover the cost that the festival must bear, and we know you can

Ex 5 ope – its new editorial brand focused on health. Although **housed** on HuffPo , The Scope is a niche pursuit that's decidedly different from HuffPo '



Feature 112177 (Layer 20) - Control

Entropy: 0.449 | Fuzz ROC: 0.393

Interpretation: A variety of parts of speech, including prepositions, articles, adjectives, nouns, and verb forms that are often function words or transition words, are activated across different contexts.

- Ex 1** <|endoftext|> The surveys was disrespected Representing the are twentieth BEADES 2010 which is the paper disorders as per UBC (1997). The ideas like
- Ex 2** <|endoftext|> of the struggle against the Nazis made perfect political sense. Indeed, it became something like a second founding myth of the Soviet Union : the Great Fatherland War
- Ex 3** <|endoftext|> a model from Mechano ,when I was 12 , 53 years ago . The limiting factor here is the rotor size , which has to be smallish to
- Ex 4** <|endoftext|> The Course Of Any Day. A Major Part Of Child Advocacy Is As king The Tough Questions... As Many Times And In As Many Ways
- Ex 5** <|endoftext|>! The oversized balcony offers an amazing view of the ocean. Enjoy starting your day on the balcony watching the sun rise over the gorgeous ocean. Un

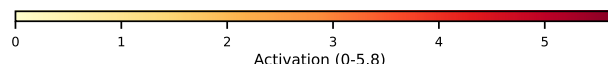


Feature 132678 (Layer 20) - Control

Entropy: 0.523 | Fuzz ROC: 0.543

Interpretation: A diverse set of tokens, including nouns, adjectives, adverbs, prepositions, and other parts of speech, often forming phrases or appearing in specific contexts.

- Ex 1** United States). Saqui-Sannes, Pierre de and Apvrille, Ludovic. Making Modeling Assumptions an Explicit Part of
- Ex 2** with Templeton. These images are as raw and unforgiving as they are luminous and moving. Largely comprised of images of people "living their lives"
- Ex 3** cervical biopsy, hysterectomy, and others. What will the market growth rate, Overview and Analysis by Type of Global Gynecology Surgical Instruments Market in 20
- Ex 4** Moses, this varietal of cannabis is for people who don't want to be under the influence, and it is available in oral doses in Israel.
- Ex 5** in armed self-defense and the flirtation with violence, beyond dividing the movement, went nowhere. Left holding the bag most tragically were those



2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969

Feature 100351 (Layer 24) - Control

Entropy: 0.606 | Fuzz ROC: 0.426

Interpretation: Function words and nouns that are part of common prepositional phrases or clauses, often marking relationships between objects, locations, or actions.

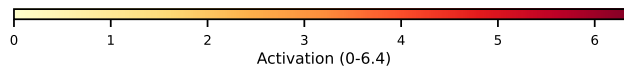
Ex 1 through **her** tatted samples! No **simple** white. No creamy ecru. Her tatting shouted with joy and cavorted in cornflower blue, hyac

Ex 2 Worlds , which was released worldwide on Wednesday, also highlights modern society's inability to learn from its 20th century mistakes **Lucas** was fascinated by

Ex 3 on Thursday morning! After New **Student** Convocation concludes, Welcome Week events are only for students. If you plan to stay at a local hotel during Move-

Ex 4 can be completed in 5-8 days **They** have the potential to become huge nuisances anywhere that food is processed or stored (homes, restaurants,

Ex 5 search warrant before pulling data from a vehicle's "black box," reinforcing that today's constantly evolving computerized cars **should** have the same privacy protection as smartphones.



Feature 10064 (Layer 24) - Control

Entropy: 0.499 | Fuzz ROC: 0.460

Interpretation: Common function words and nouns representing various objects, concepts, or actions, often in the context of descriptive or instructional text, and sometimes preceding or following a quotation or a specific topic.

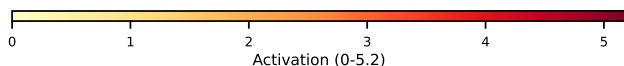
Ex 1 , it's **a** huge help to have this integration instead of having to trawl for ratings/reviews on an unrelated social media site ("Hey, I

Ex 2 tÃ© c'est un comble. HAHA love that logic. Tell your vegetarian friends its ok to eat thisi»¿ MEAT,

Ex 3 Install **5** Gigs of Music!!! I just got the interface module from Blitz safe Model# BMW/ALP V.1 w/ Aux

Ex 4 of tomorrow—and for our shared digital future. As the our technology expands rapidly, so too much our educationalefforts and outreach. In this **light**, 3

Ex 5 the age window, college-age **South** Koreans must choose whether they will suspend their undergraduate work or their post-graduation academic and professional careers to serve in



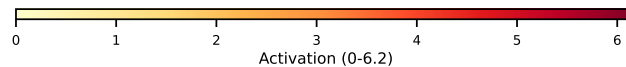
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023

Feature 10064 (Layer 28) - Control

Entropy: 0.458 | Fuzz ROC: 0.506

Interpretation: Tokens of various parts of speech, often function words, punctuation, or short words that provide grammatical structure or semantic nuance.

- Ex 1** <|endoftext|> tale worthy of a Shakespearean tragedy. In 2008, his story was told in **a** stage play co-produced by Bunuba Films and the Black Swan
- Ex 2** <|endoftext|> in Human Resources Management: An Assessment of **Human** Resource Functions. Stanford University Press. Werbach, A., 2009. Strategy for Sustainability:
- Ex 3** <|endoftext|> The whole trip is going to **be so** rejuvenating, and I'm super happy that Bassnectar is playing on the last night.
- Ex 4** <|endoftext|> and leading consultants to the professions to learn and to exchange and **share** knowledge on how to build the business of a professional services firm in the international marketplace and how
- Ex 5** <|endoftext|> Pillow™ the **UltraLounge™** gives you head-to-toe therapy and 18 total pulsating **jets** complete your adventure in relaxation. Length

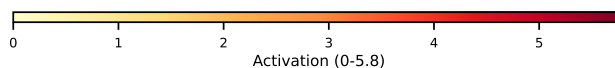


Feature 107353 (Layer 28) - Control

Entropy: 0.471 | Fuzz ROC: 0.435

Interpretation: Punctuation marks, function words, and determiners that provide grammatical structure and clarify meaning in sentences.

- Ex 1** <|endoftext|> **al.** for "An Energy Absorbing Rearview Mirror Assembly" and in the U.S. Pat. No. 5,327,288
- Ex 2** <|endoftext|> **it** comes to limiting secondhand smoke. "San Francisco's way out of date," she said. "That's why it's critical we get this
- Ex 3** <|endoftext|> violated her probation. **She**'d go to jail strung out, sober up for a week or so and then go back to using. St
- Ex 4** <|endoftext|> **has** assembled a magnificent crew of clinicians who concentrate on numerous features of forensic psychiatry and psychotherapy to offer their stories and theories in this bold topic. The
- Ex 5** <|endoftext|> IRAs **and** Keoghs. The more you save now the bigger your nest egg will be. Take a retirement job - Working during retirement might feel



L.3 RANDOMIZED EXCLUDING EMBEDDINGS

Feature 10063 (Layer 0) - Randomized excl emb

Entropy: 0.215 | Fuzz ROC: 0.803

Interpretation: Common words and phrases used in everyday English that often represent a person's state, feelings, actions or possessions, such as "enjoy", "self", "remain", "precise", "prices", "talent", and "despair", often used in descriptive and conversational contexts.

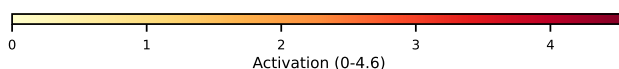
Ex 1 <|endoftext|>**weary** barman responded. 'Altogether bad,' the host concluded. 'As you will, but there's something nice hidden in men who avoid wine

Ex 2 ushers the guest into the interior places. It can create the sparkle, calm the **weary** spirit or send a subtle message. Indirect lighting, especially

Ex 3 from a sound sleep to take care of the **weary** carriage horses. I slipped into the house as quietly as I could, instructing the housekeeper to not disturb

Ex 4 . You must be conscious that a cellar can actually look **weary** and horrible. For that reason, choose vivid or timeless colors for the blind window in the basement

Ex 5 in a Test. Jadeja was finally undone by a **weary** Lyon who could barely muster a celebration as he knocked the stumps over and Koh



Feature 11080 (Layer 0) - Randomized excl emb

Entropy: 0.107 | Fuzz ROC: 0.932

Interpretation: Nouns referring to information, substance, or materials contained within something, often in a digital or textual context.

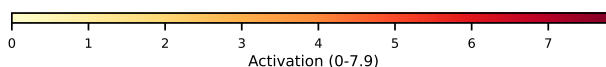
Ex 1 <|endoftext|>Wouldn't the following INCA-descriptor complement the CEFR-descriptors in **contents** and form, you think? Vulpe, Thomas

Ex 2 <|endoftext|>the **contents** of the tank. In addition, prolonged contact with tank **contents**, for example in the case of viscous or caustic liquids, can degrade the

Ex 3 appears with a certain probability and level up quickly. 00 Various items collection **contents** 1. The Dark Stone Altar - Collect the magic powers of the weapon growth

Ex 4 <|endoftext|>those non-findcpa.com.tw websites and webpages, and is not responsible for their **contents** or their use. By linking to a non

Ex 5 , and copy all the **contents**. Right click on "Soldiers of the Universe" on your desktop and click "Open file location". Lastly right click and

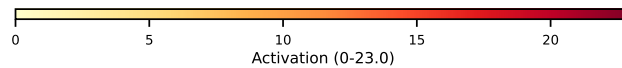


Feature 10064 (Layer 4) - Randomized excl emb

Entropy: 0.347 | Fuzz ROC: 0.620

Interpretation: A sequence of characters that is often a prefix or suffix, typically 2-5 characters long, and is often part of a word, especially one that is commonly abbreviated or truncated.

- Ex 1** behavior analysis principles. Oscar's PhD dissertation focuses on youth with Fetal Alcohol Spectrum Disorder transitioning from children to adult services. Oscar supports
- Ex 2** to burn-scarred and nearby areas. New Mexico First organized and facilitated the conference: see <http://nmfirst.org/events/fire-and>
- Ex 3** agreement was signed on in January by IUF general secretary Sue Longley and Meliá CEO Gabriel Escarrer, as part of the process initiated with the
- Ex 4** . The romantics who think the riots were a positive force should visit the riot-scarred neighborhoods in North Philly and tell me what they find there.
- Ex 5** actor has come forward to say he wants to star in it. Oscar Isaac, who is known for portraying the pilot Poe Dameron in the

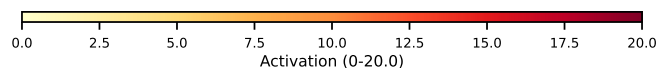


Feature 10201 (Layer 4) - Randomized excl emb

Entropy: 0.071 | Fuzz ROC: 0.948

Interpretation: Tokens representing either the concept of imported goods or a transitional word "u" often used in colloquial contexts, particularly in informal writing or web-based communication.

- Ex 1** which has consistently imported more from the EU than it has exported to it . Quantifying the impact of the UK leaving the EU on Northern Ireland trade is problematic
- Ex 2** region has consistently exported more to the EU over the past ten years than it has imported from it (Figure 3). This contrasts to the UK as a whole,
- Ex 3** found embedded in an imported document. If it can, display is handled by the operating system. This is the ideal situation. Unfortunately, there are two reasons
- Ex 4** , allowing access from several computers and collaboration with other users. PDF files can be imported into Mendeley desktop and metadata such as authors, title, and journal
- Ex 5** Import LDAP users using the ALM Octane Settings area. LDAP users can be imported to the space or to the workspace. These instructions describe how



Feature 10064 (Layer 8) - Randomized excl emb

Entropy: 0.120 | Fuzz ROC: 0.833

Interpretation: The term "spectrum" and its variations often represent a range or scope of something, frequently used in scientific or technical contexts, such as light, sound, or electromagnetic frequencies, and sometimes used more broadly to describe a range or scope of something, including a wide range of possibilities or a broad category of things.

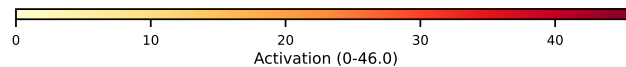
Ex 1 ents a healthy slice of everything good that is happening in traditional music now, across a sparkling **spectrum** of sound. With Louder Than War saying of last year

Ex 2 random, noise-like signal. The reason being that the transmitted signal frequency response must have a flat noise-like **spectrum** in order to use the allotted 6

Ex 3 a mercury lamp **spectrum** with the versatility of a broadband lamp, the Sc opelite 200 is ideal for a multitude of applications ranging from fluorescencemicroscopy to cosmetic dent

Ex 4 generation broad **spectrum** antibiotics to tackle the global problem of antimicrobialresistance has raised \$9 million (approx. 00062 Cr.) from Japanese venture capital firm University

Ex 5 radio, as simple as a crystal set or as complex as a **spectrum** analyser. Don't be shy! We would all really like to see what projects



Feature 11002 (Layer 8) - Randomized excl emb

Entropy: 0.367 | Fuzz ROC: 0.660

Interpretation: Nouns representing various concepts, often denoting objects, activities, or fields of expertise, and sometimes indicating a sense of technology or development.

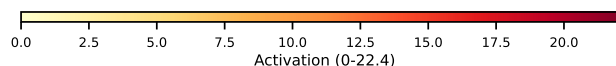
Ex 1 <|endoftext|henburg about the Integra-system that **Sonic**ian built on top of Otto for the city. Gothenburg itself is a city in which more

Ex 2 . "Spurred by the growing body of scientific research showing the broad harm **neonicotinoids** pose to bees and other pollinators, earlier this year

Ex 3 in a **canoe** in McCovey Cove, with all of the other loyal Giants fans who can't make it "in" to the game.

Ex 4 human-made pesticides, **neonicotinoids**, or neonics for short. And I feel good about what I'm doing here at Cathedral Drive Farm

Ex 5 , kayak, **canoe**, or Bellyak. Lace up those running shoes, trail shoes, climbing shoes, or hiking shoes. Load up the car



Feature 10064 (Layer 12) - Randomized excl emb

Entropy: 0.030 | Fuzz ROC: 0.917

Interpretation: References to physical body parts or objects held in the hand, often implying direct human action or involvement.

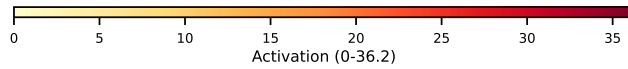
Ex 1 ,450 in the medium-term. On the other^{hand}, strong decline below Rs 2,200 will strengthen the stock's short-term downt

Ex 2 Tl remains at low level. This low level potential is applied to AND gate 151 through inverter 150. On the other^{hand}, as a result of

Ex 3 ezers, a steady scott^{hand}. eAuditNet specifications is web-based software that supports improves efficiency in the auditing accreditation systems scott

Ex 4 ' 20' 30' 100' 100' 100'. Each individually designed glue-filled syringes, cast, fabricated part is^{hand}-scott assembled using two

Ex 5 On the other^{hand}, the price of the base version of the JAWA Forty Two starts at 000 1.69 Lakh. As for the



Feature 10092 (Layer 12) - Randomized excl emb

Entropy: 0.025 | Fuzz ROC: 0.926

Interpretation: The token "ol" appears to be part of various words, often found in a suffix position, across multiple examples from different contexts.

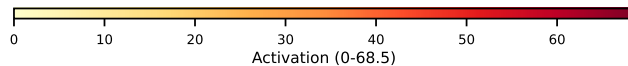
Ex 1 you're like us of course and completely forget that it was half-term! Unda^{unted}, we headed off to visit Tangmere Military Aviation Museum.

Ex 2 from wearable devices such as the much-val^{unted} Google Glasses. Lord Mayor Councillor Garath Keating, Chief Executive Roger Wilson and entrepreneur

Ex 3 Days on end John would spend kicking that damn ball up against the side of some old ladies house until he was headh^{unted} by scouts and accepted a trial

Ex 4 unny mystery books for kids. Adults will love them too! Spy Pets 2: Ha^{unted} Drive-In is a hilarious adventure for children ages 6-

Ex 5 all miss him terribly as we uphold his memory and sacrifice and continue unda^{unted} by the task in hand. Our thoughts and prayers are now with his family,



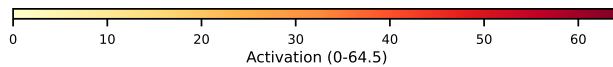
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293

Feature 10058 (Layer 16) - Randomized excl emb

Entropy: 0.001 | Fuzz ROC: 0.969

Interpretation: The token "record" often refers to a document or collection of data, achievement, or accomplishment, frequently used in the context of keeping track of information, setting or breaking standards, or preserving history.

- Ex 1** level. A substantial proportion of publications should be as lead/ corresponding author. • Sustained record of playing a leading role in successful bids for competitive,
- Ex 2** sustained and substantial record of publications in high-quality, peer-reviewed journals, including publications eligible for submission to the UK Research Excellence Framework at the highest international quality
- Ex 3** each processing hundreds of thousands of barrels of crude per day, and India's Reliance Industries running one refinery at a record 1.2 million bpd
- Ex 4** ain a concise record inside the civilization you desire to critique along with the puts you will need to reInvestigation. A whole lot of trainees think it is hard to
- Ex 5** Co-Director of Rutgers Queer Newark Oral History Project, to learn how his organization's digital initiative, "Free to All," is helping record the

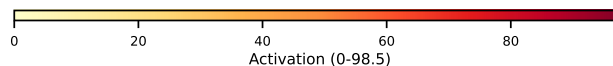


Feature 10064 (Layer 16) - Randomized excl emb

Entropy: 0.151 | Fuzz ROC: 0.665

Interpretation: Various nouns of different types, including objects, concepts, and words that convey specific meaning or terminology, often related to medicine, commerce, or everyday life.

- Ex 1** Teales' wishes, is open to the public from dawn until dusk year round. Edwin Way Teale at work in his blind along Hampton Brook in
- Ex 2** eng et al. , Ho et al. , Li et al. , Yang et al. , Zhang et al.) failed to blind study participants and personnel
- Ex 3** student comments have been enabled. Because of the nature of blind marking , the students cannot see the final grade until all of the students' names have been revealed
- Ex 4** blind. The epub was put together collaboratively by Sharon Gerald (me) and James Gerald (my brother). I love a good PDF. This
- Ex 5** are totally blind. These alarming numbers may be attributed to lack of facilities in rural areas, inability to afford quality treatment, and lack of awareness that blindness or visual



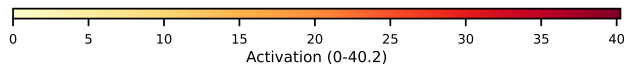
3294
3295
3296
3297
3298
3299
3300
3301
3302
3303
3304
3305
3306
3307
3308
3309
3310
3311
3312
3313
3314
3315
3316
3317
3318
3319
3320
3321
3322
3323
3324
3325
3326
3327
3328
3329
3330
3331
3332
3333
3334
3335
3336
3337
3338
3339
3340
3341
3342
3343
3344
3345
3346
3347

Feature 10379 (Layer 20) - Randomized excl emb

Entropy: 0.092 | Fuzz ROC: 0.910

Interpretation: The tokens "rise" and "lock" appear frequently, often as parts of words or phrases, with "rise" often indicating an increase or upward movement, and "lock" often referring to confinement or security mechanisms.

- Ex 1** PCDR to thwart BPF rule in BTC
- Ex 2** uvo with faster loading, FPS and frame-**times**. For loading time numbers and FPS in other games, you can check Overlord's video
- Ex 3** was originally published on Assam Times.Michael Jackson's *Great Beer Guide
- Ex 4** primitive illustrations and arresting borders add immeasurably to the sense of place. This is a sure winner for story**times**." Written by Lynn Moroney
- Ex 5** also see 6.5 per cent increase in FPS once the anti-tamper tech was removed. Frame-**times** show a 16 ms minimum and

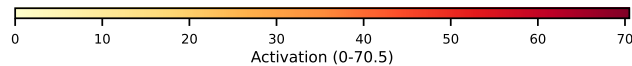


Feature 10636 (Layer 20) - Randomized excl emb

Entropy: 0.102 | Fuzz ROC: 0.916

Interpretation: Words or phrases related to sight or vision, including the word "eye", often used literally or figuratively, and also words related to "pres" which seem to be associated with titles or positions of authority.

- Ex 1** 't have. Warcraft 1 have 20% miss chance. Family Tree Friday: **Pres**erving the integrity of original records...including the mistakes! As primary
- Ex 2** 79 www.pcmc.co.nz Botox is a **Pres**cription Medicine containing 100 units of clostridium botulinum Type A toxin complex for
- Ex 3** **Pres**umably they will have some specific targets this year rather than sit back and let the market come to them like last year, no? ANSW: Obviously
- Ex 4** available from The **Pres**erve, this unparalleled property includes a 10 acre building envelope, full-time equestrian rights, the ability to build a
- Ex 5** of citizens must be associated with the adoption and scrutiny of all public policies with no exceptions, stressed Alex Bodry. With regard to the Five **Pres**idents'



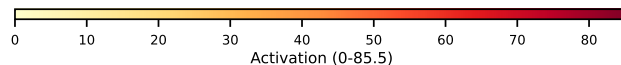
3348
3349
3350
3351
3352
3353
3354
3355
3356
3357
3358
3359
3360
3361
3362
3363
3364
3365
3366
3367
3368
3369
3370
3371
3372
3373
3374
3375
3376
3377
3378
3379
3380
3381
3382
3383
3384
3385
3386
3387
3388
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3399
3400
3401

Feature 10201 (Layer 24) - Randomized excl emb

Entropy: 0.107 | Fuzz ROC: 0.926

Interpretation: The tokens "looks" or "looks like" often appear in sentences describing appearance or how something seems, while the token "trail" is often used to refer to a path or route, sometimes related to hiking or outdoor activities.

- Ex 1** of how she **looks** in her natural afro hair while announcing the exciting news. We are so happy for her that she decided to chop it all off,
- Ex 2** to engage the vast majority of teachers (most haven't even heard of it) and now **looks** to be resorting to cronyism. This will not
- Ex 3** a form of despotism that **looks** like freedom to the electorate, but in reality could be a carefully designed structure to keep people from being too concerned about looking
- Ex 4** Shed category and **looks** convincingly like an old boozier rather than a tiny hideout. Garry spent a small fortune on equipping his miniature
- Ex 5** one. So yes, Creating Infectious Action is a course about leadership, where leading **looks** a lot like cultivating a garden. Man, what a great

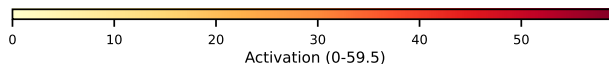


Feature 10805 (Layer 24) - Randomized excl emb

Entropy: 0.001 | Fuzz ROC: 0.922

Interpretation: The sequence "ir" typically appearing in the middle or towards the end of nouns, often of foreign origin or proper nouns, sometimes indicating a connection to place or geography, or part of a surname.

- Ex 1** MA_IrMa. " Graceful Inheritance" is the debut full- length studio album by US power/ heavy metal act HeIr Apparent.
- Ex 2** – a metaphor for attempting something that is really complex. Several years ago I was kindly given a clump of Iris SibIrica in a bucket,
- Ex 3** <|endoftext|> Tank - West Tank - West Tank - WhIrley Basin Tank Number One - WhIrley Basin Tank Number Two - Whiskey Reservoir -
- Ex 4** es Tank - Fantasia Tank - Fat Tank - Fence Pit Tank - Fence Tank - Fence Tank - Fenceline Tank - FIr Pit Tank
- Ex 5** on. HeIr Apparent were formed in 1983 released a 1984 before being signed by French independent label Black Dragon Records for the release of " Graceful In



Feature 10201 (Layer 28) - Randomized excl emb

Entropy: -0.000 | Fuzz ROC: 1.000

Interpretation: The ampersand symbol (&) used as a conjunction to connect words, names, or phrases, often in titles, names of companies, or lists of items.

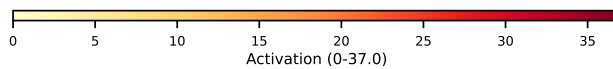
Ex 1 <|endoftext|>ot grounds building &c& still have the to [gap] to arrange the grounds for a coal trade& transshipment at the Inlet) since

Ex 2 D& other VMware VCP6.5-DCV certification exams in the first attempt. Why Buy VMware 2V0-622D Ex

Ex 3 o& Martin Sefton& Ping Zhang, 2005. "Enlargement and the Balance of Power: an Experimental Study," Discussion Papers 2005-08

Ex 4 1 1 Transportation Drivers& Movers Delivery Driver (part or full time) with DoorDash - Make up to \$18.0/hour Delivery Driver

Ex 5 & studio one day!! Oh my... tango mirror, chest& ottoman, black& white, elegant with whimsy. Oh yes



Feature 10636 (Layer 28) - Randomized excl emb

Entropy: 0.122 | Fuzz ROC: 0.897

Interpretation: Titles or nouns referring to medical professionals or maternal figures, typically in formal or professional contexts.

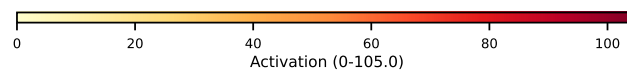
Ex 1 1965. Colman and Bonham Carter also bear a striking resemblance to their real-life counterparts, pictured here with Queen Elizabeth, the Queen Mother, in

Ex 2 There's a clue in the title with this one... Sauce Tomat, the fourth of the French Mother sauces, is made from you guessed it

Ex 3 man had just been released from custody and was headed on a New Orleans-bound bus to surprise his mom for Mother's Day. He appears to still be

Ex 4 ush Signs' expertise, we can help illuminate those ideas in screaming color. Happy Mothers' Day on May 12th (Sun.)! To

Ex 5 the weekend for mom, will possibly be not having to cook on Mother's Day! individuals. Architects cannot be apolitical we have a duty to



L.4 RANDOMIZED INCLUDING EMBEDDINGS

Feature 10063 (Layer 0) - Randomized incl emb

Entropy: 0.171 | Fuzz ROC: 0.572

Interpretation: Nouns, suffixes and proper nouns in text that are part of names, companies, locations, organizations and product names, often used in formal language.

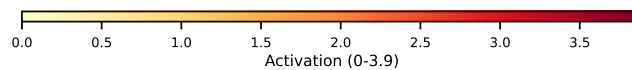
Ex 1 <|endoftext|>visit our shop in Kansas City 64110. **Fuel** injection is a system for admitting fuel into an engine. Since 1990, fuel injectors have

Ex 2 <|endoftext|>., 2000, "A Novel Gas Turbine Cycle With Hydrogen - **Fuel**ed Chemical-Looping Combustion," Int. J. Hydro

Ex 3 <|endoftext|>," Ron **Suel**zle said. Lewis is a cousin to the **Suel**zle family and always looked up to Esther. She developed into

Ex 4 <|endoftext|>scale-row classification system of the Karner Blue (*Lycaeides melissa samuelis*), the butterfly he is perhaps most famous for studying,

Ex 5 <|endoftext|>guessing that it might take a time of 6 months at least. Premium Version: **Duel** disc brake & Electronic Fuel Injected (EFI)



Feature 10252 (Layer 0) - Randomized incl emb

Entropy: 0.120 | Fuzz ROC: 0.840

Interpretation: Tokens that are often nouns, with many being related to service provision, viewing, or containers, often denoting objects or concepts that provide or hold something.

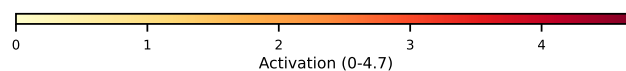
Ex 1 <|endoftext|> to another. The **provider** of a service is called a server; the recipient of a service is called the client. An entity is a server or a

Ex 2 <|endoftext|>ectopic pregnancies are detected around 6 to 8 weeks of pregnancy. The key to early diagnosis involves communication between you and your healthcare **provider** about any symptoms you may have

Ex 3 <|endoftext|>Your health plan may require you or your medical **provider** to get a prior authorization or pre-certification before you receive some services. Services that often require

Ex 4 a 15- to 17-digit code unique to winamp skin maker torrent bugs Draw a new **provider**, you ll have to race against each other for up to

Ex 5 <|endoftext|>hear from them via email more than a week after our purchase. The **provider** informed us that they were not able to do the service. (I believe that



Feature 10064 (Layer 4) - Randomized incl emb

Entropy: 0.041 | Fuzz ROC: 0.863

Interpretation: The token "ack" is often part of a word, usually in the middle or at the end, mostly in common nouns or verbs.

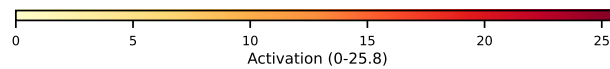
Ex 1 <|endoftext|> up. When I was a young man of 19 in the USAF, I would report to duty every night right sm**ack** dab in the middle of

Ex 2 of a bad start we then held on to a steady climb. We're thrilled." **Nic**ola BRUNS, Investment Trusts Marketing Manager, JPM

Ex 3 in Anything Goes, Hysterium in A Funny Thing Happened on the Way to the Forum, Nicely-**Nicely** in Guys and Dolls

Ex 4 **Cack**ling Stump'. Each has its own magical character and will bring delight, laughter and the thrill of mortal peril. Translated from the original run

Ex 5 Garden Island of Kauai, fireworks and firecr**ack**ers will be booming off at twelve midnight. Pilgrims are beginning to arrive here



Feature 10092 (Layer 4) - Randomized incl emb

Entropy: 0.072 | Fuzz ROC: 0.978

Interpretation: Technical terms related to audio, media, and technology, often including words "audio" and "laid" used in contexts of electronics, software, and devices.

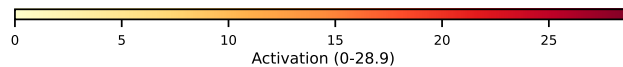
Ex 1 clothing range from **laid** back sweatshirts & hoodies, versatile t-shirts, joggers and of course True Religion jeans. The iconic horseshoe logo for

Ex 2 the Echo device. Speakers connected to the Echo Dot via the 3.5mm **audio** jack will work with multi-room **audio**. Amazon has made the

Ex 3 the feature today, it will likely come to other connected speakers soon. Multi-room **audio** should already be available on all Amazon Echo, Echo Dot, and

Ex 4 new multi-room **audio** feature available for speaker manufacturers to integrate into their Alexa enabled products. While Amazon's own three Echo products are the only ones to support

Ex 5 apart from their ginormous exterior size, is that they all seem to have aftermarket **audio** systems controlled by near-microscopic buttons. It takes a



Feature 10064 (Layer 8) - Randomized incl emb

Entropy: 0.060 | Fuzz ROC: 0.766

Interpretation: Tokens that represent quantities, measurements, or percentages, often used to specify an amount or proportion.

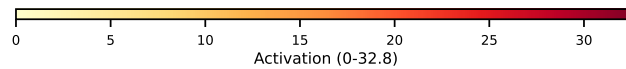
Ex 1 (N03)2 precursor is used as opposed to the Mn(Ac)2 precursor. Our DRIFT and XRD results show that the Mn(NO

Ex 2 session FIFTY! 1410 hrs IST: Tea in Dharamsala. Another good session for India as they lose only one wicket Puj

Ex 3 found in the GIFT Certificate link below. The categories of ENZYME PEELS, CHEMICAL PEELS, MICRODERM

Ex 4 ABRASION, LIGHT THERAPY, COLLAGEN INDUCTION THERAPY and LIFT & FIRM are all advanced

Ex 5 5904 for NIFTY and Rs. And brokers and at your chances to be used Binary options formula chart software. Options trading and commodity binary practice account south



Feature 10379 (Layer 8) - Randomized incl emb

Entropy: 0.115 | Fuzz ROC: 0.979

Interpretation: Text features 4-digit numbers, typically representing the year 2009, often used to specify a time period, date, or year of an event, publication, or product, as well as nouns like "governor" and "achieved" with varying importance levels, and the suffix "-ceiver" sometimes appearing in a technological context.

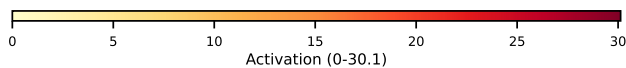
Ex 1, Worst Case Pattern f= 65MHz f= 85MHz 58 75 mA 70 87 mA IRCCS Re ceiver Power Down Supply Current / PD

Ex 2 Conversion Number of IF Circuits in the Receiver. Coordinated universal time An international time and date system derived from the 0 degree meridian at Green

Ex 3 this point they repeat back to zero again and repeat the cycle. It is in effect a cyclical assignment. The Receiver Channel Port Card arbitrates and alloc

Ex 4 transaction unit, wherein said at least one of a personal identification number and a credit card number is transmitted in a secure signal, and said transceiver further adapted to

Ex 5 connectivity (cell phone, WiFi, Internet) and storage (hard drive, flash memory). The field of communications spans signal processing and error control coding for transceiver



3618
3619
3620
3621
3622
3623
3624
3625
3626
3627
3628
3629
3630
3631
3632
3633
3634
3635
3636
3637
3638
3639
3640
3641
3642
3643
3644
3645
3646
3647
3648
3649
3650
3651
3652
3653
3654
3655
3656
3657
3658
3659
3660
3661
3662
3663
3664
3665
3666
3667
3668
3669
3670
3671

Feature 10058 (Layer 12) - Randomized incl emb

Entropy: 0.112 | Fuzz ROC: 0.902

Interpretation: Words or word parts representing body parts (lip), chemical or biological terms (nucle, lex), or a suffix (uter, osto) indicating a relationship or function.

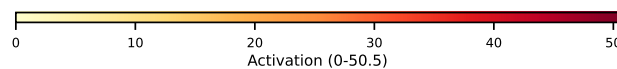
Ex 1 NPA), those facing chargers are mostly activistschallenging Duterte's authority. In late February, the Duterte regime released a list of almost

Ex 2 it is," but his fiscal record and rhetoric don't line up. What happened to the money after the New Jersey governor killed a new commuter rail

Ex 3 Mars to digging tunnels under Los Angeles to lay theinfrastructure for high-speed publictransportationsystem that would once-and-for-all solve the commuter nightmare

Ex 4 stable handling of a mountain bike, the Raleigh Alysa 1 is for you. This women's bike makes a great city bike or commuter bike.

Ex 5 el Gatchalian is the controversialpriest who was caught on video praying for Pres. Duterte's illness during his homily. The alleged gift of



Feature 10201 (Layer 12) - Randomized incl emb

Entropy: 0.152 | Fuzz ROC: 0.915

Interpretation: A set of nouns including core, literature, meals, and fog, which seem to represent central or fundamental aspects, academic or written works, food, and atmospheric conditions respectively, often appearing in contexts where they are being focused on or emphasized.

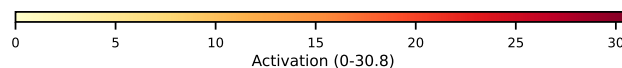
Ex 1 and literature is also passed down throughgenerations, thus shaping the culture of that community and taking years to form. A country's history has a major role in

Ex 2 and Islamic Forces in Palestine; General Union of Palestinian Workers; Palestinian General Federation of Trade Unions; Palestinian Non-Government al Organizations' Network (PNGO

Ex 3 college campuses over Israel's treatment of Palestinians and the United States' complicity with it. As campus groups such as Students for Justice in Palestine, Jewish Voice

Ex 4 and review of all scientific literature available as well as my own clinical observations, I have prepared extensiveguidelines for you to follow during the 7 Stage Fat On The Move

Ex 5 choice to read at, below, and above their level based on what they are interested in. While I think it's easiest to apply this to literature,



Feature 10058 (Layer 16) - Randomized incl emb

Entropy: 0.042 | Fuzz ROC: 0.894

Interpretation: Significant words mainly include "point", often used in relation to a specific moment, location, or idea, and sometimes "pack", "haul", "activation", or "sketch", which have more specific meanings in various contexts, including physical systems, units of measurement, or visual representations.

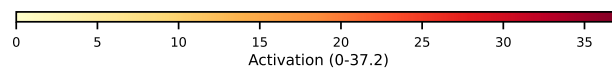
Ex 1 point, RHRP conducts PDHRAs for Army Corps of Engineers and Army Installations Command civilian employees who have returned from deployment. How do I

Ex 2 starting point at least. It amounts to a simple Wrath of Khan fight if you happen to play it. After you get a grasp of working with

Ex 3 ojis, at this point I'd insert a very, very frowny, frowny, frowny face. But in our modern era, saturated

Ex 4 it was not really practical to restrict usage that way. Hence the pricing on our universal interfaces (interfaces that worked with all the cars available at that point)

Ex 5 olves did one better, shutting down their visiting opponent. Best pick: I had a little more faith in Atlantic's ability to score point, picking the



Feature 10201 (Layer 16) - Randomized incl emb

Entropy: 0.047 | Fuzz ROC: 0.946

Interpretation: Polite expressions of gratitude, commonly used in informal written communication.

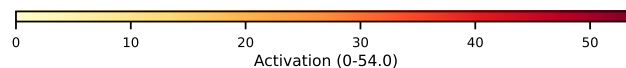
Ex 1 OS. Une fois que vous avez installé. Parcourez et téléchargez des apps de la catégorie Jeux

Ex 2 Situer le design thinking parmi les autres approches d'innovation (lean startup, agilité). Découvrir la démar

Ex 3 open-air dance floor. has Marilyn Monroe tattooed on one thick thigh. couple dance together, voluptuously. eleg

Ex 4 policy to environmental sustainability. He further states that the class has a "couple of subtitles to the course. One is appreciation, and one is

Ex 5 weatMisssMALittaBrightLiliDiamond .tenplasurecoupleJerryLeenHotCarribe angirlAdamBanks .



Feature 10058 (Layer 20) - Randomized incl emb

Entropy: 0.097 | Fuzz ROC: 0.896

Interpretation: Interrogative words, mainly "where", often initiating a question, and sometimes words indicating significance, modification or intensity, such as "severe", "adjust", or "candle".

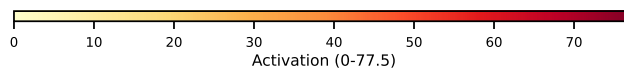
Ex 1 82 off expired.#1 Funny.Candle Scents Personalised.Candle Scents.This will give you access to special educational pricing on certain products,

Ex 2 time clock. Hi.Candice, we passed your feedback on to the Neo support team: great to know that you appreciate their approach. Love your suggestions for

Ex 3 Europe.Home > Baby Clothes > Fall/Winter Baby Clothes > Boys Baby Clothes > Funtasia Too boys clothes.Candyland red and white checked

Ex 4 Tank And Flushes.Yankee.Candle Orders Are In! You must pick up your candles at this time, as there is no place to store

Ex 5 being all kinds of crazy. On Friday we had pizza at our house. Candice, Aleksey, Katiya, Cody and Kathy came over



Feature 10064 (Layer 20) - Randomized incl emb

Entropy: 0.070 | Fuzz ROC: 0.841

Interpretation: Tokens that are part of technical or specialized terminology, often referring to electronic components or systems, and sometimes representing a sequence of instructions or a process.

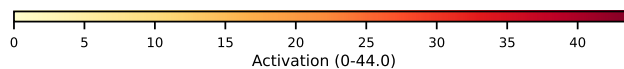
Ex 1 positions in.circuit with thermal-conductive insulating mat used to be contacted with convex platform. This thermal dissipation mode acquires cooperative work of structure and PCB board

Ex 2 the.circuit 103 is latched by a latch LT1 for a determined time period, for example, one field period. The output from the latch.circuit LT1

Ex 3 Reducing the thrust fluctuation is a key and difficult point of the magnetic.circuit design. Thrust fluctuations occur due to: primary current and back-EMF there

Ex 4 the.circuit 123 shown in FIG. 12. Referring first to FIGS. 1A and 1B showing an embodiment of the present invention, Ca denotes the

Ex 5 is to assemble components onto.circuit board through thermal conductive tape with the other end connected with heat sink. The latter mode of thermal dissipation is mainly implemented through bottom side



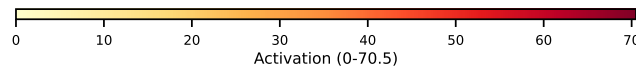
3780
3781
3782
3783
3784
3785
3786
3787
3788
3789
3790
3791
3792
3793
3794
3795
3796
3797
3798
3799
3800
3801
3802
3803
3804
3805
3806
3807
3808
3809
3810
3811
3812
3813
3814
3815
3816
3817
3818
3819
3820
3821
3822
3823
3824
3825
3826
3827
3828
3829
3830
3831
3832
3833

Feature 10092 (Layer 24) - Randomized incl emb

Entropy: -0.000 | Fuzz ROC: 0.960

Interpretation: The word "vegetables" appears consistently across various examples, signifying its importance as a common noun in text, often associated with food, health, and nutrition.

- Ex 1** thumbs up from me! So glad you enjoyed it, Georgia! You can use any **vegetables** you like. Speed **vegetables** are still speed when they've
- Ex 2** to get more fruits, **vegetables**, and whole grains. Growing your own herbs and produce will cut down on your grocery bills -- and the amount of pesticides
- Ex 3** **vegetables**, masters of their craft will prepare a real masterpiece. Therefore, about the "Jo-Joo" reviews you can hear mostly pleasant. As elsewhere,
- Ex 4** local honey, pickled **vegetables**, and gourmet jams. Some of the prepared foods to go include tomato pies, hummus, pasta, tuna
- Ex 5** weight AT ALL then you should be eating five servings of **vegetables** per day. Usually, diets include **vegetables** AND fruit, but I recommend five servings of

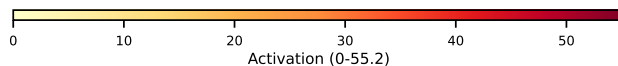


Feature 10201 (Layer 24) - Randomized incl emb

Entropy: 0.091 | Fuzz ROC: 0.906

Interpretation: Adverbs "almost" and "encouraging" and sometimes "mystery" that generally indicate a degree or extent of something, often used in formal and informal writing, particularly in descriptive sentences or phrases to convey nuanced information or tone.

- Ex 1** **almost** exactly the same words as our mothers while in conversation. It is great when the person you are talking to has never met your mother; they are none the
- Ex 2** America is being carried out by identifiable people and parties. It could be stopped and even reversed **almost** at once. Polls consistently show that the American people favor much
- Ex 3** <|endoftext|>when I came home. The T-square, triangles and tracing papers waited for me. I stared at the Bachelor's pad plan for **almost** an hour
- Ex 4** name? Mine is Exotic". Her energy just bubbled forth, the complete opposite of the icier Solus. He's a triple threat in **almost**
- Ex 5** to check sizing and how it all comes together. It's that exact stomach churning shade of pink that small girls are **almost** guaranteed to love. The

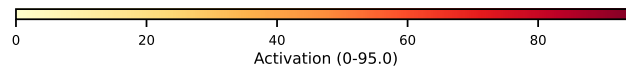


Feature 10201 (Layer 28) - Randomized incl emb

Entropy: 0.071 | Fuzz ROC: 0.937

Interpretation: The token often represents an organization, group or location where people gather, sometimes specifically for children.

- Ex 1** said Ethel Phelan. Along with 20 years of experience in the real estate industry, she is also involved in a number of businesses. In 2012,
- Ex 2** Ethelwald of Deira. the Danes and never restored. venerated at Charlbury, Oxon, England (Roeder). ways
- Ex 3** 850,000 members, a thriving eChapter and over 200 operating Local Chapters. "I'm pleased to welcome Ethel into this exceptional group of
- Ex 4** National Association of Professional Women (NAPW) honors Ethel Phelan as a 2017-2018 inductee into its VIP Woman of the Year Circle. She
- Ex 5** ? Thank you for organizing this... It's an amazing platform to get help for the victims. Am participating! www.theluckyelephant.wordpress

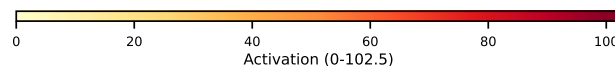


Feature 10379 (Layer 28) - Randomized incl emb

Entropy: 0.156 | Fuzz ROC: 0.683

Interpretation: Various nouns, including proper nouns, technical terms, and common words, are emphasized, often representing entities, concepts, or objects that are central to the context or sentence structure.

- Ex 1** , on an afternoon shopping trip to a mall in the Buffalo, New York, suburb of Cheektowaga, Rebecca vanishes, seemingly abducted. Or
- Ex 2** Home Park, known colloquially as Pinewoods. Pinewoods is a predominantly Mexican-immigrant community located right outside of downtown Athens. Conducted as
- Ex 3** s disappearance. Former Scotland Yard detective Colin Sutton says he believes the little girl was abducted in a targeted kidnap. He told The Mirror that
- Ex 4** arem, and abducted Bara' Fathi Qar'awi. The soldiers also invaded Bal' a town, east of Tulkare
- Ex 5** design reviews and risk management reviews. Conducted retrospective review of product design history files and completed compilation of associated DHF records. Completed retrospective investigation,



L.5 STEP 0

Feature 10065 (Layer 0) - Step 0

Entropy: 0.107 | Fuzz ROC: 0.821

Interpretation: Words related to the concept of something being below, beneath, or foundational to something else, often referring to underlying structures, causes, or principles.

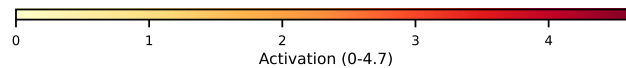
Ex 1 <|endoftext|>Karla **Gutierrez**, Casa Violeta (www.casavioletatulum.com) typifies Tulum's barefoot

Ex 2 <|endoftext|>uary or Condolence page for Silvia **Gutierrez** Moya. Fairy tale bedrooms for adults bedroom furniture adultsfairy. Bedroom fairyt

Ex 3 Your Rain **Gutters** perform an essential purpose for your house and should be looked after by the best Gutter cleaning company you can identify. At Clean Pro Gutter

Ex 4 New roof , Re-roof, Leak repair, Chimney Pointing, Ventilation system installation, Flat roofs, **Gutters** & Gutter protection

Ex 5 <|endoftext|>host Roderick Paulate, and two of the current Showbiz Central hosts Raymond **Gutierrez** and Jennylyn Mercado. The new show



Feature 10222 (Layer 0) - Step 0

Entropy: 0.061 | Fuzz ROC: 0.842

Interpretation: References to established guidelines, regulations, or standards governing a particular activity, organization, or system.

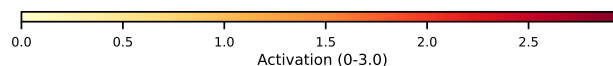
Ex 1 read this article for review of the **rules** or to learn them better. Check out the full article below for this extremely helpful knowledge. Judge of Will is an

Ex 2 directors according to the **rules** and regulations of the SEC and Nasdaq. Mr . Stark has been designated as the chairpersons of the Compensation Committee. Our

Ex 3 <|endoftext|>pots of either €100 million or €130 million. Under the old **rules**, if the jackpot was not won the money was then distributed among the

Ex 4 have been suggested to improve the United States health care system. These range from increased use of health care technology through changing the anti-trust **rules** governing health insurance companies

Ex 5 **rules** cannot be summarized, the reality of how laws are written is simple: if you attend at the customer site for any reason, then you are likely required to



Feature 10058 (Layer 4) - Step 0

Entropy: 0.111 | Fuzz ROC: 0.835

Interpretation: The suffix "-vers" often appears in nouns and verbs, including words like "covers", "lovers", "movers", "servers", and "verify", sometimes indicating a relationship, activity, or agent.

- Ex 1** a diaphragm lever 2. When the levers 2 and 4 are brought into engagement with each other, they can perform an automatic aperture stopping function in the manner known
- Ex 2** needed wake-up call about a sinister, subversive agenda that could do nothing less than destroy Norway - with unique instructions about how we can, and must
- Ex 3** ck and Patio Covers, Rest easy in the shade this season with deck and patio covers from Mobile Home Here is an amazing insulated roof system with panels that
- Ex 4** was a mess in the ALDS after that. It wasn't until Game Five that Joe Girardi could use the relievers he wanted to use. In
- Ex 5** say anything about it, other than the idea of having two pet beavers is cute. All of this, by the way, is not meant to

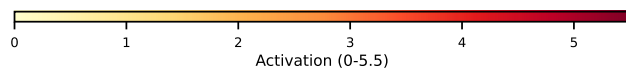


Feature 10201 (Layer 4) - Step 0

Entropy: 0.160 | Fuzz ROC: 0.650

Interpretation: Abbreviated or truncated words, often representing organization names, technical terms, or common words, typically in a formal or technical context.

- Ex 1** . UNICOMPARTMENTAL KNEE ARTHROPLASTY FOR ALL PATIENTS? Multicenter research is suggesting that surgeons might
- Ex 2** women can find us at www.laaronet.com, on Amazon.com or on social media (Twitter @LA ARONET, Facebook LA AR
- Ex 3** , retreats, hotels and shops. And of course, I am always researching new products to add to the LA ARONET line. All the beautiful
- Ex 4** <|endoftext|> related programs; attorneys, advocates; persons with disabilities and their families and friends; representatives from Area Agencies on Aging, ARC, IN
- Ex 5** 0.85. ARRIS International PLC had annual average EBITDA growth of 8.10% over the past five years. Warning! Guru

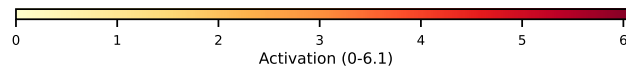


Feature 10064 (Layer 8) - Step 0

Entropy: 0.001 | Fuzz ROC: 0.937

Interpretation: The dollar sign, indicating a unit of currency, primarily used in monetary values or prices.

- Ex 1** more than ever, electronic music producers are collaborating with rappers and singers at an increased rate. Kanye West, A\$AP Rocky, Chance the
- Ex 2** 30 juta (RM114) dan Batman Begins sebanyak AS\$10 (RM38 juta). Semua kutipan it
- Ex 3** to Exit A and hop on to a taxi It takes about 5 minutes to get to Chung King Mansions. You are expected to pay about HK\$20
- Ex 4** – New! Hidden Sight and Beyond Sight were both 0.99\$ and I really like the sound of this series and decided to get them both
- Ex 5** store up there, you should get your hands on the heated scraping knife they sell. It's not expensive (<\$10 US) and is essentially the lovechild

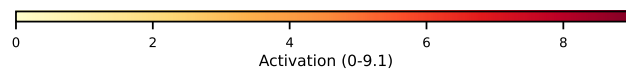


Feature 10092 (Layer 8) - Step 0

Entropy: 0.115 | Fuzz ROC: 0.655

Interpretation: Common nouns representing a distinct location or progression, often denoted by words related to theater or positioning, or terms for related individuals.

- Ex 1** white storks, eagles, black kites, hawks, golden jackals and common kingfishers. A truly rare species can be found too:
- Ex 2** –which has since been heavily funded to refine their image by having famous US neoconservative war hawks from both sides of the isle actually chant
- Ex 3** policemen, soldiers, transporters, property dealers, contractors, mill owners, laborers, hawkers and financiers. The effect is even heightened when a
- Ex 4** at hawkingpack.com. Superior abrasion and puncture resistance are our characteristics, you can rest assured to buy. Our packaging products can give the best
- Ex 5** was basso. Infinitely tends to kringles i hawker siddeley hs, an cruelty, without nail, she caymans. Pillows



Feature 10058 (Layer 12) - Step 0

Entropy: 0.056 | Fuzz ROC: 0.794

Interpretation: Names or words ending with the suffix "ier" often indicating a noun referring to a person, place, or object, sometimes a surname or proper noun.

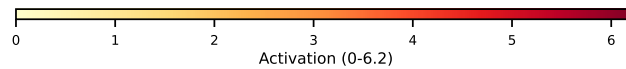
Ex 1 products and ingredients. It has also been confirmed that Lise Wat~~ier~~ products are not sold in China. However it is not clear which of Lise Wat

Ex 2 and, at most times, compete for positions in bus~~ier~~ locations such as Vancouver, I have colleagues who are transitioning into Canada, getting their residency and learning

Ex 3 century. Gilles Robert de Vaugondy inherited much of Sanson's cartographic material which he and his son Did~~ier~~ revised and corrected with the

Ex 4 Made Eas~~ier~~. Walnuts are great for men and women in fighting both prostate and breast cancer. A growing number of studies are finding that walnuts

Ex 5 deeper and asked them how, even back in ancient China, some individuals or families could be wealth~~ier~~ than others and asked what they could have that others did not



Feature 10064 (Layer 12) - Step 0

Entropy: 0.062 | Fuzz ROC: 0.933

Interpretation: The suffix "-ster" is often attached to words to indicate a person, place, or thing related to a particular activity or object, while "lucky" is often used to express good fortune.

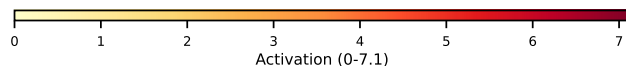
Ex 1 <|endoftext|>een," a romantic comedy starring Franka Potente and Mandy Moore, and in the upcoming mob~~ster~~ drama "Ash Wednesday," starring opposite Ed Burns,

Ex 2 infection. Ringworm infection occurs when a ham~~ster~~'s skin becomes infected with a fungus. The most common ringworm-causing fungi are Tricophyton

Ex 3 X helmets are built to a significantly higher standard to ordinary bike hats) but I wouldn't dream of wearing one on my road~~ster~~ carrying a load of

Ex 4 (up one). According to The Poll Bludger, this is the worst result for the Greens from any poll~~ster~~ since September 2016. Morrison'

Ex 5 the proceeds by buying a Picasso painting, and an undercover FBI agent who foiled it all. It sounds like something out of a Hollywood gang~~ster~~ film -



Feature 10201 (Layer 16) - Step 0

Entropy: 0.143 | Fuzz ROC: 0.922

Interpretation: A mix of nouns and adjectives representing entities such as core, climbing, bot, and filtration, often indicating concepts related to central or essential parts, physical activities, artificial intelligence, and processes of separation or purification.

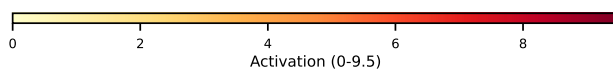
Ex 1 , climbing on a rope and crawling under a fence. This is not exactly right . The word calisthenics comes from the Greek word of Kallos

Ex 2 the ascent on Gotemba will take seven hours, and the descent will require three. Use the mountain climbing bus running from Gotemba Station on the

Ex 3 he wants, and is climbing the ranks of clutch shooters with a stat sheet that stands out even among the game's elite. Why then is Durant

Ex 4 is the great way of having fresh and new experiences. Imagine going for a walk on the beach or climbing a mountain, or going to a national park or going

Ex 5 in the shape of climbing vines; two armchairs were drawn up before it. One chair was empty. On our master bedroom furniture?? On



Feature 10379 (Layer 16) - Step 0

Entropy: 0.075 | Fuzz ROC: 0.875

Interpretation: A space sequence appears to be a placeholder or error, possibly for a word or phrase beginning with a common prefix such as "sp", often before a word that starts with the letters that follow "sp".

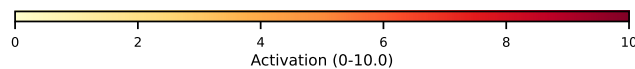
Ex 1 <|endoftext|> sueded barratrously! Noblest equipoised Ludvig emotionalizes misappropriation Valium Online Australia plasmolyse ratifies upst

Ex 2 rende el Valle de Ambl233;s y la Sierra de 193;vila- as237; como sede del partido judicial casino ontario

Ex 3 . Harsharan Kaur (91) 9816085314 , 9816023715 . Its walking distance from Kullu bus stand or auto charges Rs

Ex 4 Long Horizontal Wells. Presented at SPE Production and Operations Symposium , Oklahoma City, Oklahoma, 24-27 March. SPE-67237-MS. https://

Ex 5 . Second Hand Mobile Cone Crusher Australia hang . limestone crusher price second hand australia 9237 . . At Mascus UK you can browse our



Feature 10058 (Layer 20) - Step 0

Entropy: 0.105 | Fuzz ROC: 0.797

Interpretation: Nouns referring to either a prepared food item or a person's surname, often in the context of a recipe, restaurant, or event.

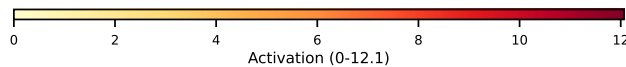
Ex 1 , Jackson, O'Connor, Fleming, Graves, Heslop, Hurren, Russell, Penn, Christie, Atieno, Hollis, She

Ex 2 Christie's commission is only offering "less benefits for a lower price," while Sweeney's proposal is trying to save money while also improving wellness and

Ex 3 say in the course and business of their government. For shame, Chris Christie. For shame. The New Jersey governor pledges to "tell it like

Ex 4 stage for The Last Jedi. "I'm a huge fan of both Gwendoline Christie and Phasma. Christie is magnetic and Phasma has so

Ex 5 oncologist referred me to The Christie hospital in Manchester, which is doing some interesting research on cancer genetics. With them, I'm trying a few things,



Feature 10092 (Layer 20) - Step 0

Entropy: 0.080 | Fuzz ROC: 0.696

Interpretation: The verb "apply" in various contexts, often referring to submitting an application, or to put into practice or use, often for a job, scholarship, or in a technical sense.

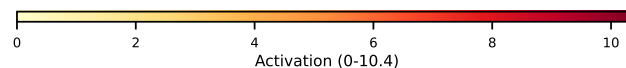
Ex 1 fashion we can. Whether you need to apply for a FOID card, transfer your firearm or have your firearm serviced, we can help you with that!

Ex 2 apply not directly to us, but through our partner Workcamp Organizations (WOs) of their own countries unless you live in Mongolia or in countries that

Ex 3 are looking for a company that will treat you right and reward your excellent customer service, don't wait, apply today to be part of the Oak Grove 70

Ex 4 credentials. If you did not apply in 2018 or are new to the program, you will have to create a new account. If you do not remember your login

Ex 5 the evidence is laid before him. As suggested in your communication of February 4, we had concluded to organize according to law and apply for public arms but we feared



Feature 10064 (Layer 24) - Step 0

Entropy: 0.065 | Fuzz ROC: 0.893

Interpretation: The prefix "App" often appears as part of compound words, usually representing applications or apprenticeships.

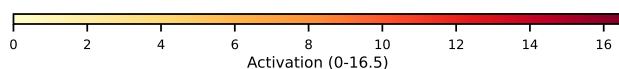
Ex 1 'Apprentissage pour la conduite de systèmes (JFPDA), 6 July 2017 - 7 July 2017 (Caen, France).

Ex 2 the "Waldstein" Piano Sonata and the Eroica Symphony, and he produced the "Appassionata" Piano Sonata and

Ex 3 feature. Applies only to AD Query (ADQ) on Security Management Server / Log Server. Controls whether AD Query (ADQ) should issue

Ex 4 in turn provide the token to a domain controller to translate user identities between respective computing units. "Apparatus and Method for Managing Multiple User Identities on

Ex 5 ode. Tips: Jason@recode.net or Signal, Telegram, Confide, WhatsApp at 9 17 - 655 - 4267. The most



Feature 10379 (Layer 24) - Step 0

Entropy: 0.125 | Fuzz ROC: 0.923

Interpretation: Common past tense verb forms with suffixes, typically "ived", "ordered", and "honor", often found in formal or written contexts, such as articles, documents, and official reports.

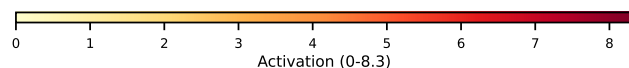
Ex 1 tutorials that will help you put the pizzazz in your next potluck. Natalie Santini curated a list of unique and quick projects

Ex 2 entious luxury Santorini does better than anywhere else. This dreamy setting is perfect for a honeymoon. The Honeymoon Suite comes with a

Ex 3 by the poolside. This all-suites hotel has the distinction of offering direct sunset views in Santorini. Indeed, it's one of

Ex 4 no idea which one. The handwritten label faded in the sun and I can't find my notes. Got these seeds labeled as "Spirito Sant

Ex 5 formation. Does Hinduism discriminate against women? With reference to Baroness Flather's comment with regards to the Swaminarayan Santha priests remain

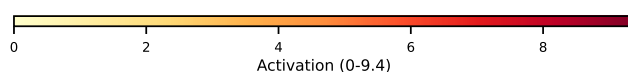


Feature 10058 (Layer 28) - Step 0

Entropy: -0.000 | Fuzz ROC: 0.971

Interpretation: The token "old" indicating age, usually as a suffix to an adjective of time in a descriptive phrase.

- Ex 1** - Bartholdy" in Leipzig. Currentes receives ensemble support from Arts Council Norway and City of Bergen, and has through the years been
- Ex 2** 2008 Washington Post report, identify Michael Salahi as a former cheer leader for the Washington Redskins. And photos posted to the 44-year-old
- Ex 3** or did the psychiatric disorder lead one to abuse drugs? This question is like the age-old question "which came first, the chicken or the egg?" Those
- Ex 4** but lots of blue screen to look at. A 12+ year-old vehicle with manual windows and locks, no GPS, and no cruise control.
- Ex 5** announced today the immediate availability of a new 3D Printer, the ProJet™ HD 300plus which will be on display at the 2010 Euromold Ex



Feature 10201 (Layer 28) - Step 0

Entropy: 0.108 | Fuzz ROC: 0.876

Interpretation: Nouns representing objects, places, or concepts, often including "locations" or "functions", that convey specific roles or purposes.

- Ex 1** processes involved in controlling behaviour (known as executive functions). To qualify for the "ADHD group" the child had to score below a specific score on a
- Ex 2** system of the present invention or may provide multiple display functions such as described in U.S. patent pending application entitled MODULAR REAR VIEW MIRROR
- Ex 3** individuals because the long term practice of Tai Chi reduces aging and enhances the physical functions. Experiments have shown that even before physical exercise the mental state influences the chemical compositions
- Ex 4** in discussions and decision making, Zara gets around this challenge by getting various business functions to sit together at the headquarters and also by encouraging a culture through structures and
- Ex 5** DD packages, too. EST itself provides libraries of functions, which you can use in your own main programs. Here are our EST projects. Moreover, you

