KRONECKER FACTORIZATION IMPROVES EFFICIENCY AND INTERPRETABILITY OF SPARSE AUTOENCODERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Sparse Autoencoders (SAEs) have demonstrated significant promise in interpreting the hidden states of language models by decomposing them into interpretable latent directions. However, training and interpreting SAEs at scale remains challenging, especially when large dictionary sizes are used. While decoders can leverage sparse-aware kernels for efficiency, encoders still require computationally intensive linear operations with large output dimensions. To address this, we propose Kron-SAE – a novel architecture that factorizes the latent representation via Kronecker product decomposition, drastically reducing memory and computational overhead. Furthermore, we introduce mAND, a differentiable activation function approximating the binary AND operation, which improves interpretability and performance in our factorized framework.

1 Introduction

Interpreting large language models and their embeddings in particular remains a central challenge for transparency and controllability in AI systems (Elhage et al., 2023; Heap et al., 2025). Sparse autoencoders (SAEs) have emerged as powerful tools for uncovering human-interpretable features within neural activations by enforcing activation sparsity to induce discrete-style dictionaries (Elhage et al., 2023; Gao et al., 2025; Cunningham et al., 2024). These dictionaries facilitate circuit-level semantic analysis (Marks et al., 2025) and concept discovery, enabling fine-grained probing of model internals (Elhage et al., 2023; Cunningham et al., 2024).

However, naively scaling SAEs to the widths demanded by modern transformers leads to prohibitive compute costs, limiting their applicability to large-scale interpretation experiments. *Gated* SAEs (Rajamanoharan et al., 2024a) address this by learning continuous sparsity masks via lightweight gating networks, achieving a Pareto improvement on the reconstruction–sparsity trade-off. *Switch* SAEs (Mudide et al., 2025) leverage conditional computation by routing activations among smaller expert SAEs, reducing computation by activating only a subset of experts per input. Gao et al. (2025) propose an acceleration of Topk SAE that utilizes an optimized kernel based on efficient sparse–dense matrix multiplication. Encoder remains unoptimized: it still performs a dense projection into the full dictionary, incurring high computational cost and limiting scalability.

Another limitation of SAEs is the absence of structure within learned latents in its classical design. While Mudide et al. (2025) address this via expert subnetworks, Bussmann et al. (2025) imposes feature hierarchy via loss function and improves the interpretability of SAE latents.

In this paper we address both these directions and introduce **KronSAE**. By decomposing the latent space into head-wise Kronecker factors and including differentiable logical AND-like gating mechanism, KronSAE reduces both parameters and compute overhead while preserving reconstruction fidelity and improves feature interpretability.

This work makes three primary contributions:

1. We identify the encoder projection as one of the principal scalability bottlenecks in sparse autoencoders, demonstrating that targeted encoder optimizations can significantly improve computational performance while maintaining reconstruction quality.

- 2. We propose **KronSAE**, a Kronecker-factorised sparse autoencoder equipped with the novel mAND activation function. The design reduces encoder cost and is compatible with existing sparse decoder kernels, loss function designs and mixture of experts architecture.
- 3. We show on multiple language models that KronSAE decreases feature absorption and yields more interpretable latents under fixed compute.

2 RELATED WORK

 Sparse Autoencoders. Early work demonstrated that SAEs can uncover human-interpretable directions in deep models (Cunningham et al., 2024; Templeton et al., 2024), but extending them to the large latent sizes $(F \gg d)$ for modern language models with large hidden dimension d is expensive: each forward pass still requires a dense $\mathcal{O}(F\,d)$ encoder projection. Most prior optimization efforts focus on the decoder side: for example, Gao et al. (2025) introduce a fused sparse–dense Topk kernel that reduces wall-clock time and memory traffic, while Rajamanoharan et al. (2024a) decouple activation selection from magnitude prediction to improve the ℓ_0 -MSE trade-off. Separately, Rajamanoharan et al. (2024b) propose Jumprellu, a nonlinearity designed to mitigate shrinkage of large activations in sparse decoders.

Conditional Computation. These schemes avoid instantiating a full dictionary per token by routing inputs to a subset of expert SAEs via a lightweight gating network (Mudide et al., 2025) by employing the Mixture-of-Experts ideas (Shazeer et al., 2017), but still incur a dense per-expert encoder projection, leaving the encoder as the primary bottleneck.

Weight Factorization and Logical Activation Functions. Outside of sparse decoders and Mixture-of-Experts, tensor-factorization methods have been used to compress large weight matrices in language models (Edalati et al., 2021; Wang et al., 2023). In parallel, differentiable logic activations were introduced to approximate Boolean operators in a smooth manner (Lowe et al., 2021). Our method synthesizes these lines of work: we embed a differentiable AND-like gate into a Kronecker-factorized encoder to build compositional features while preserving end-to-end differentiability.

3 Method

Preliminaries. Let $\mathbf{x} \in \mathbb{R}^d$ denote an activation vector drawn from a pretrained transformer. A conventional TopK SAE (Gao et al., 2025) produces a reconstruction $\hat{\mathbf{x}}$ of \mathbf{x} via

$$\mathbf{f} = \text{TopK}(W_{\text{enc}}\mathbf{x} + \mathbf{b_{enc}}), \qquad \hat{\mathbf{x}} = W_{\text{dec}}\mathbf{f} + \mathbf{b_{dec}},$$
 (1)

where $W_{\text{enc}} \in \mathbb{R}^{F \times d}$ and $W_{\text{dec}} \in \mathbb{R}^{d \times F}$ are dense matrices and $\mathbf{f} \in \mathbb{R}^F$ is a sparse vector retaining only the K largest activations. The encoder cost therefore scales as $\mathcal{O}(Fd)$ per token.

KronSAE. Our method reduces the encoder's computational cost while also enforcing compositional structure of the latents. We decompose the latent space into h independent heads, and each head k is parameterised by the composition of two thin matrices $P^k \in \mathbb{R}^{m \times d}$ (composition base) and $Q^k \in \mathbb{R}^{n \times d}$ (composition extension), with dimensions $m \le n \ll d$ and F = h m n. The pre-latents

$$\mathbf{p}^k = \text{ReLU}(P^k \mathbf{x})$$
 and $\mathbf{q}^k = \text{ReLU}(Q^k \mathbf{x}),$ (2)

acting as the elements from which compositional features would be built, are combined through an element-wise interaction kernel independently in each head:

$$z_{i,j}^k := \text{mAND}(p_i^k, q_j^k) := \begin{cases} \sqrt{p_i^k q_j^k}, & p_i^k > 0 \text{ and } q_j^k > 0, \\ 0, & \text{otherwise,} \end{cases}$$
 (3)

where $\mathbf{z}^k \in \mathbb{R}^{m \times n}$; it is then flattened in a row-major order to a vector equivalent to the element-wise square root of the Kronecker product of \mathbf{p}^k and \mathbf{q}^k . Assuming that $\text{vec}(\cdot)$ is in row-order, we get

$$\operatorname{vec}(\mathbf{z}^k) = \sqrt{\operatorname{vec}(\mathbf{p}^k(\mathbf{q}^k)^\top)} = \sqrt{\mathbf{p}^k \otimes \mathbf{q}^k}$$
 (4)

Concatenating heads and applying TopK yields *post-latents* $\mathbf{f} \in \mathbb{R}^F$. This decomposition induces a hierarchical structure, where \mathbf{q} and \mathbf{p} latents can be interpreted as fundamental, lower-level features; the mAND kernel smoothly approximates Boolean AND gate, ensuring non-zero activation only if *both* inputs are positive, while square root preserves the activation magnitude for stable reconstruction by preventing activation value explosion when both operands are positive.

The encoder cost per token drops from $\mathcal{O}(Fd)$ to $\mathcal{O}(h(m+n)d)$ (see Appendix A.2). KronSAE thus reduces FLOPs and parameter count without routing overhead, and is orthogonal to existing sparse decoder kernels (Gao et al., 2025) and thus can be combined with them for end-to-end speed-ups.

4 EXPERIMENTS

We train SAEs on the residual streams of Qwen-2.5-1.5B-Base (Yang et al., 2024), Pythia-1.4B-deduped (Biderman et al., 2023), and Gemma-2-2B (Team et al., 2024) language models. Activations are collected on FineWeb-Edu, a filtered subset of educational web pages from the FineWeb corpus (Penedo et al., 2024). We measure reconstruction quality via explained variance (EV),

$$EV = 1 - \frac{Var(\mathbf{x} - \hat{\mathbf{x}})}{Var(\mathbf{x})},$$

so that 1.0 is optimal, and use automated interpretability pipeline (Bills et al., 2023; Paulo et al., 2024) and SAE Bench (Karvonen et al., 2025) to evaluate properties of SAE features. We aim for the needs of resource-constrained mechanistic interpretability research where efficiency and interpretability are in favor rather than top reconstruction performance, and 100M-820M token budgets are widely adopted (Bussmann et al., 2025; Kharlapenko et al., 2025; Heap et al., 2025; Mudide et al., 2025; Karvonen et al., 2025), so we choose 125M, 500M and 1B token budgets for the experiments.

Our experiments (see detailed setup in Appendices A and D) address three questions:

- 1. Does KronSAE maintain EV comparable to TopK SAE under fixed compute?
- 2. Which design choices (nonlinearity, (m, n, h)) drive EV improvements?
- 3. How do these choices affect properties and interpretability of learned latents?

4.1 ABLATIONS

We employ the iso-FLOPs setup: for each KronSAE variant of dictionary size F we allocate the same amount of FLOPs as was spent for the training of TopK SAE for token budget T and same F.

Reconstruction performance. As indicated on Figure 1, KronSAE achieves on-par performance with TopK given lower number of trainable parameters. The performance gap narrows when increasing the dictionary size, which indicate the potential scalability of our method for large dictionaries.

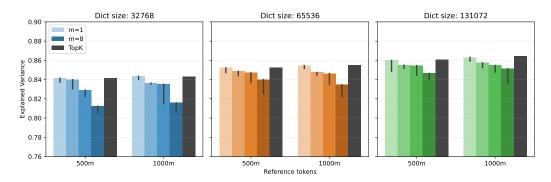


Figure 1: Maximum performance for KronSAE vs. TopK SAE on Qwen-1.5B for different dictionary sizes F and budgets in iso-FLOP setting. KronSAE with lower number of parameters achieves performance comparable to the baseline, and the gap narrows with larger dictionary size.

Decomposition hyperparameters. We systematically vary the number of heads h and the per-head base dimension m (with n=F/(mh)) under the iso-FLOPs setup. From the Figure 2, we conclude that lower m and higher h consistently yields higher reconstruction quality, due to flexibility of pre-latents - as we show in section 5.3, they must be either expressive or fine-grained enough (low m, n and high h) to efficiently represent the semantics. See also Appendices A.3 and A.5.

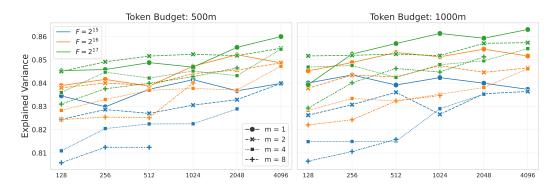


Figure 2: Dependency of EV on head count h and base dimension m under 500M and 1B token budgets in iso-FLOPs setup. Higher h and smaller m yield improved reconstruction quality because of higher expressivity of pre-latents to encode semantics and increasing trainable parameters.

Composition activation. To isolate the impact of our mAND operator, we compare it to two simpler interaction kernels: (i) the element-wise product of ReLUs, $\operatorname{ReLU}(u) \cdot \operatorname{ReLU}(v)$, and (ii) the raw product $u \cdot v$. As reported in Table 1, under a 125M training budget, the mAND variant achieves the highest explained variance. More description of the mAND is provided in the Appendix C. See also our experiments where we replace TopK with JumpReLU in Appendix A.4.

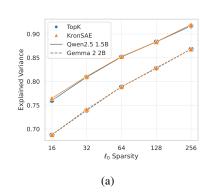
Dictionary size	m	n	h	Activation	Explained Variance
32768	2	4	4096	$ \begin{vmatrix} \text{mAND}(u, v) \\ \text{ReLU}(u) \cdot \text{ReLU}(v) \\ u \cdot v \end{vmatrix} $	0.8336 0.8267 0.8237
	4	8	1024	$ \begin{vmatrix} \text{mAND}(u, v) \\ \text{ReLU}(u) \cdot \text{ReLU}(v) \\ u \cdot v \end{vmatrix} $	0.8220 0.8191 0.8143
65536	2	4	8192	$ \begin{vmatrix} \text{mAND}(u, v) \\ \text{ReLU}(u) \cdot \text{ReLU}(v) \\ u \cdot v \end{vmatrix} $	0.8445 0.8328 0.8297
65536	4	8	2048	$ \begin{vmatrix} \text{mAND}(u, v) \\ \text{ReLU}(u) \cdot \text{ReLU}(v) \\ u \cdot v \end{vmatrix} $	0.8350 0.8297 0.8251

Table 1: Performance of different composition activations under a budget of 125M tokens.

Sparsity and depth-wise position. Figure 3a compares KronSAE and TopK SAE across a range of ℓ_0 sparsity settings on the 14th layer of Qwen-2.5-1.5B and the 12th layer of Gemma-2-2B; Figure 3b evaluates performance across layers in Qwen-2.5-1.5B. In every case, KronSAE matches the reconstruction quality of the TopK baseline, demonstrating that our Kronecker-factorized encoder maintains its performance regardless of sparsity level or depth.

4.2 ABSORPTION

The notorious challenge in SAE interpretability is *feature absorption*, where one learned feature becomes a strict subset of another and consequently fails to activate on instances that satisfy the



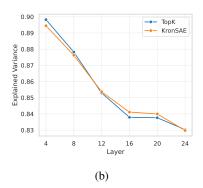


Figure 3: (a) EV versus sparsity level ℓ_0 for KronSAE and TopK SAE on the 14th layer of Qwen-2.5-1.5B and the 12th layer of Gemma-2-2B under iso-FLOPs constraints. (b) EV across layers of Qwen-2.5-1.5B, demonstrating that KronSAE matches TopK performance regardless of depth.

broader concept but not its superset representation (e.g. a "starts with L" feature is entirely subsumed by a "Lion" feature) (Chanin et al., 2024).

Figure 4 reports three absorption metrics measured via SAEBench (Karvonen et al., 2025) across sparsity levels $\ell_0 \in \{16, 32, 64, 128, 256\}$: (1) the *mean absorption fraction*, measuring the proportion of features that are partially absorbed; (2) the *mean full-absorption score*, quantifying complete subsumption events; and (3) the *mean number of feature splits*, indicating how often a single conceptual feature fragments into multiple activations. Across all ℓ_0 , KronSAE variants consistently reduce first two scores relative to the TopK SAE baseline, while maintaining a similar rate of feature splits.

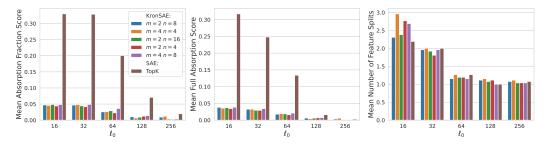


Figure 4: Feature absorption metrics on Qwen-2.5-1.5B. KronSAE configurations (various m, n) exhibit lower mean absorption fractions and full-absorption scores across different ℓ_0 .

We attribute KronSAE's improved disentanglement to two complementary design choices:

- 1. **Smooth** mAND **activation.** By emitting nonzero output only when both pre-latents are positive, we introduce a differentiable AND gate that prevents a broadly polysemantic primitive from entirely subsuming a more specific one. Consequently, composite post-latents fire mainly in the intersection of their constituent concepts, encouraging each pre-latent to specialize on a single semantic mode rather than inherit its "parent" activation region.
- 2. **Head-wise Cartesian decomposition.** Dividing the latent space into h independent subspaces (each with its own $m \times n$ grid of primitive interactions) ensures that specialized concepts (such as "elephant") are confined to a single head and cannot fully absorb more general concepts (such as "starts with E") in another.

Together, these mechanisms produce more monosemantic features, as we verify in the section 5.3, simplifying downstream applications. Notably, the mean number of feature splits remains comparable to the TopK baseline, as decomposition alone does not inherently alter the fragmentation of individual primitives (see section 5.2).

This result is consistent with observations made for Matryoshka SAE (Bussmann et al., 2025), which also impose structure of the learned latents significantly reducing the absorption score, although using

not architectural innovations but training methodology. We also validate the result on Gemma and Pythia models and observe the same picture, see the Appendix A.6.

5 ANALYSIS

In this section we examine the properties of the latents in KronSAE compared to TopK architecture.

The design of our arhictecture was also inspired by the observation that many features within a single SAE correlate with each other. By imposing the compositional structure via our encoder design, we force post-latents within each head to correlate; we expect that this would allow KronSAE to move correlated features into the same head.

By examining the toy examples with manufactured correlations in data, we show that KronSAE captures these correlations better than TopK. Then we show that KronSAE trained on language indeed moves correlated features within a single head, indicated by higher correlation within head. After that, we show that KronSAE pre-latents interactions are closely resemble the logic AND gate, and its post-latents are notably more interpretable than TopK latents.

5.1 TOY MODEL OF CORRELATION

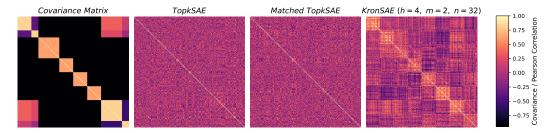


Figure 5: We generate data with covariance matrix that consist of blocks with different sizes on diagonal and off diagonal (left panel). We then examine the decoder-weight covariance $W_{dec} \cdot W_{dec}^{\top}$ to assess feature-embedding correlations. Panels (middle panels) show that a TopK SAE recovers these correlation structures only weakly, even after optimal atom matching. In contrast, KronSAE (right panel) more accurately reveals the original block patterns.

To evaluate how well different sparse autoencoder architectures recover known correlation patterns, we construct a controlled experiment using a synthetic, block-structured covariance model. Input vectors $\mathbf{x}_{\mathrm{sparse}} \in \mathbb{R}^{256}$ sampled under a heavy-tailed Bernoulli distribution with probability p=0.875 that value will be zero. We perform a Cholesky decomposition $S=LL^{\top}$ on the covariance matrix S and set $\bar{\mathbf{x}}_{\mathrm{sparse}} = L\,\mathbf{x}_{\mathrm{sparse}}$, so that $\bar{\mathbf{x}}_{\mathrm{sparse}}$ exhibits the desired structure.

We train autoencoder (AE) to reconstruct $x_{\rm sparse}$ following the (Elhage et al., 2022):

$$\hat{\mathbf{x}} = \text{ReLU}(W^{\top}W \cdot \bar{\mathbf{x}}_{\text{sparse}} + \mathbf{b}).$$
 (5)

We collect hidden states $(W \cdot \bar{\mathbf{x}}_{\mathrm{sparse}})$ from AE and then train TopK SAE and our proposed KronSAE with |F| = 256 and topk = 4 to reconstruct it. After training, we extract the decoder weight matrices W_{dec} from each model and compute the covariance $C_{\mathrm{dec}} = W_{\mathrm{dec}}W_{\mathrm{dec}}^{\top}$. To compare TopK SAE embeddings to the AE reference, we match atoms by minimal Euclidean distance, ensuring fair alignment before analysis. Result is shown in Figure 5.

To quantify how closely each model's feature correlations mirror the ground-truth structure S, we employ the RV coefficient, defined as $RV(S,C)=\mathrm{trace}(SC)/\sqrt{\mathrm{trace}(S^2)\,\mathrm{trace}(C^2)}$ and assess its significance via a permutation test. In our experiments, KronSAE consistently achieves RV=0.358 with p-value = 0.0002 while TopK SAE which achieved RV=0.038 with p-value = 0.31 does not show any structure at all. Also, we try to find the nearest pattern in TopK SAE that matches its feature embeddings with learned features in AE. This setup has better score RV=0.080 with p-value = 0.001, but still much less than KronSAE.

These results indicates that our compositional encoder more faithfully reconstructs the original feature relation. See also additional experiments in section B, where we isolate the architectural bias from the results and verify that KronSAE identifies correlations when they are actually presented in the data.

5.2 CORRELATIONS IN SAES TRAINED ON LANGUAGE

To examine the correlation structure of features learned in our SAE, we have calculated the correlations on 5k texts from the training dataset. For each feature we calculate the mean correlation with features within its head and with all other features, and compare the randomly initialized KronSAE with m=4, n=4 with the trained one. To isolate the effect of our initialization procedure, we initialize the weights of SAE from the uniform distribution. As shown in Figure 6, correlations are indeed significantly higher within a single head and higher than for random SAE, which suggest that our choice to impose the correlated structure in SAE latents works as intended.

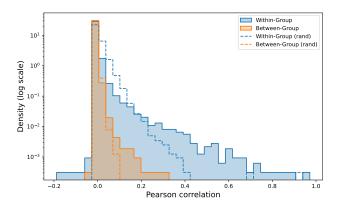


Figure 6: Correlations between features in KronSAE with m=4, n=4 within a head and with features from other heads. Our design induces higher correlations within a group, which also gets stronger after training, although SAE have also learned correlated features from different heads.

5.3 Analysis of Learned Features

In this section we compare KronSAE and TopK SAE in terms of interpretability and feature properties, and we analyze the properties of groups in KronSAE. For this, we choose the 14th layer of Qwen2.5-1.5B and a dictionary size of 32k features, of which the first 3072 were selected. KronSAE was chosen with m=4, n=4. We run for 24M tokens total to collect data. Our interpretation pipeline follows the common methodology: LLM interprets the activation patterns (Bills et al., 2023) and we evaluate obtained interpretations using the *detection* score and the *fuzzing* score Paulo et al. (2024).

For each selected feature, among the standard mean activation value and frequency, we calculate two additional metrics. Low values of *token entropy* suggest that feature activates more frequently on small number of tokens, thus it is token-specific; high value of *multitoken ratio* indicates that feature tends to activate multiple times in a single sentence. We have observed that both these metrics have notable negative correlation with the final interpretability scores and therefore they provide useful signal to assess the potential score without calculating it.

For more details on the data collection and interpretation pipeline, see Appendix D. For additional analysis of properties of learned features, see Appendix E.

SAE properties and encoding mechanism. We observe that the features learned by KronSAE are more specific, indicated by lower values of the computed metrics and higher interpretability scores, as shown in Figure 7. Since post-latents are significantly more interpretable than corresponding prelatents, we hypothesize the hidden mechanism for encoding and retrieval of the required semantics.

By examining activating examples and interpretations of latents, we observe that pre-latents may carry multiple distinct and identifiable modes of activation, such as composition base element 3 in head 23 shown in Table 2, and be very abstract compared to resulting post-latents. Polysemanticity of

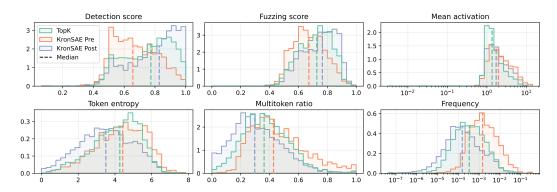


Figure 7: Distribution of properties for TopK SAE and KronSAE (m=4, n=4) with 32k dictionary size trained on Qwen2.5-1.5B. Pre and Post suffixes denote pre- and post- latents, and y-axis indicate density. Our SAE achieves better interpretability scores by learning specialized feature groups, indicated by lower activation frequency and lower variance in activated tokens.

pre-latents is expected to be a consequence of reduced "working" number of encoder latents, since we decompose the full dictionary size and reduce the encoder capacity.

Thus, we hypothesize that the encoding of specific semantics in our SAE may be done via magnitude, which we validate by examining the activation examples. For the above mentioned pre-latent, the "comparison" part is encoded in the top 75% quantile, while the "spiritual" part is mostly met in the top 25% quantile, and the "geographical" part is mainly encoded in the interquartile range. We also consider but do not investigate the possibility that it may depend on the context, e.g. when the model uses the same linear direction to encode different concepts when different texts are passed to it.

Semantic retrieval and interpretable interactions. Heads usually contain a groups of semantically related pre-latents, e.g. in head 136 there are three base elements and one extension covering numbers and ordinality, two extension elements related to geographical and spatial matters, one question-related base and one growth-related extension. Interestingly, most post-latents for this head have higher interpretability score than both its parent pre-latents, which is unusual.

The retrieval happens primarily via the mechanism closely resembling the logical AND circuit, where some pre-latent works as the bearer of multiple semantics, and the corresponding pre-latent (base or extension) works as specifier. An illustrative example is shown in Table 2: we see that the base contains three detectable sub-semantics, and each extension then retrieves the particular semantics.

Other types of interaction may occur, such as appearance of completely new semantics, for example composition between base 3 and extension 1 in Table 2 where medical terminology arises and could not be interpreted as simple intersection between two pre-latents semantics. Another example is a case of head 3 where base 3 has sub-semantics related to technical instruments and extension 2 have semantics related to the possession and necessity, and their combination gives the therapy and treatment semantics which looks more like addition than intersection.

It is a frequent case that post-latent inherit semantics of only one parent, or the impact of another parent is not detectable, which usually happens if parent has a very broad interpretation and low score. However, it requires more sophisticated techniques to properly identify the fine-grained structure of interactions than just looking at the resulting latent descriptions, so we leave it to further work. Despite this, the AND-like gate is a very common behavior. See more examples in Appendix G.

Geometry of post-latents. Each post-latent vector has a vector representation in the residual stream represented by the corresponding column in $W_{\rm dec}$, which is the approximation of overcomplete basis vectors we search for when training SAEs. We had not observed any notable differences in feature geometry between TopK and our SAEs, except that our architectural design leads to clustering so that post-latents produced by same head, base or a extension elements are grouped in a tight cluster, and the geometry is dependent on hyperparameters h, m, n we choose, which is expected and may be useful for further applications such as steering. See more details in Appendix E.

Component	Interpretation	Score			
Base 3	Suffix "-like" for comparative descriptors, directional terms indicating geographical regions, and concepts related to spiritual or metaphysical dimensions				
Extension ele	ments and their compositions with base 3				
Extension 0	<i>Interpretation:</i> Comparative expressions involving "than" and "as" in contrastive or proportional relationships.	0.87			
	<i>Composition:</i> Similarity or analogy through the suffix "-like" across diverse contexts.	0.89			
Extension 1	<i>Interpretation:</i> Specific terms, names, and abbreviations that are contextually salient and uniquely identifiable.	0.66			
	Composition: Medical terminology related to steroids, hormones, and their derivatives.	0.84			
Extension 2	<i>Interpretation:</i> Spiritual concepts and the conjunction "as" in varied syntactic roles.	0.80			
	Composition: Abstract concepts tied to spirituality, consciousness, and metaphysical essence.	0.93			
Extension 3	<i>Interpretation:</i> Directional and regional descriptors indicating geographical locations or cultural contexts.	0.84			
	Composition: Directional terms indicating geographical or regional divisions.	0.91			

Table 2: Interactions between composition base element 3 in head 23 and all extension elements in that head. Interaction happens in a way that closely resembles the Boolean AND operation: base pre-latent is polysemous, and the composition post-latent is the intersection, i.e. logical AND between parent pre-latents. See details in Section 5.3.

6 CONCLUSION AND FUTURE WORK

We introduce **KronSAE**, a sparse autoencoder architecture design that combines head-wise Kronecker factorization of latent space with a approximation of logical AND via mAND nonlinearity. Our approach allows to efficiently train interpretable and compositional SAE, especially in settings with limited compute budget or training data, while maintaining reconstruction fidelity and yielding more interpretable features by utilizing their correlations. Our analysis links these gains to the complementary effects of compositional latent structure and logical AND-style interactions, offering a new lens on how sparsity and factorization can synergise in representation learning.

Limitations. KronSAE introduces tradeoff between interpretability, efficiency and reconstruction performance, and due to reduced number of trainable parameters it is expected to lag behind TopK SAE at large budgets. Our evaluation is limited to mid-sized transformer models and moderate dictionary sizes; however, the main bottleneck there might be not the SAE itself, but the infrastracture required to handle these setups and the model inference.

Future Work. We identify three directions for extending this work: (i) *Transcoding*. Treat transcoders (Dunefsky et al., 2024) as implicit routers of information and investigate alternative logical gating functions (e.g. XOR or composite gates) to improve interpretability and circuit analysis. (ii) *Crosscoding*. Generalize KronSAE to a crosscoder setting (Lindsey et al., 2024) uncover interpretable, cross-level compositionality via logic operations. (iii) *Dynamic Composition*. Explore learnable tuning of both the number of attention heads and their dimensionality, enabling fine-grained decomposition into groups of correlated features at varying scales.

ETHICS STATEMENT

While interpretability research has dual-use potential, our method operates within the ethical boundaries of the underlying models and aims to advance responsible AI development through better model understanding. We analyze activations from publicly available language models (Qwen-2.5-1.5B, Pythia-1.4B, and Gemma-2-2B) gathered on FineWeb-Edu datasets, which excludes the unreported harmful content. We declare no conflicts of interest and maintain transparency about limitations, including potential artifacts from LLM-based interpretation as noted in Appendices D.3 and I.

REPRODUCIBILITY STATEMENT

We have taken several measures to ensure the reproducibility of our results. We use publicly available models (Qwen, Gemma, Pythia families) and training dataset (FineWeb-Edu) in our experiments. Section 4 and Appendix A provide detailed description of SAE training procedure and hyperparameter configuration. Our complete implementation is available in the supplementary materials, containing the training code, interpretation pipeline and analysis of the results. Appendix H includes simplified implementation of KronSAE that might be easily integrated into existing training codebases, while Appendix D details the interpretability analysis methodology with precise evaluation protocols.

REFERENCES

- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models across training and scaling, 2023. URL https://arxiv.org/abs/2304.01373.
- Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. Language models can explain neurons in language models, 2023. URL https://openai.com/index/language-models-can-explain-neurons-in-language-models/.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, et al. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2, 2023.
- Bart Bussmann, Noa Nabeshima, Adam Karvonen, and Neel Nanda. Learning multi-level features with matryoshka sparse autoencoders. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=m25T5rAy43.
- David Chanin, James Wilken-Smith, Tomáš Dulka, Hardik Bhatnagar, and Joseph Bloom. A is for absorption: Studying feature splitting and absorption in sparse autoencoders, 2024. URL https://arxiv.org/abs/2409.14507.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. In *International Conference on Learning Representations (ICLR)*, 2024. URL https://openreview.net/forum?id=F76bwRSLeK.
- Jacob Dunefsky, Philippe Chlenski, and Neel Nanda. Transcoders find interpretable LLM feature circuits. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=J6zHcScAo0.
- Ali Edalati, Marzieh Tahaei, Ahmad Rashid, Vahid Partovi Nia, James J. Clark, and Mehdi Rezagholizadeh. Kronecker decomposition for gpt compression, 2021. URL https://arxiv.org/abs/2110.08152.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henigan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022. URL https://transformer-circuits.pub/2022/toy_model/index.html.

- Nelson Elhage et al. Decomposing language models with dictionary learning. *Trans-former Circuits*, 2023. URL https://transformer-circuits.pub/2023/monosemantic-features.
 - Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=tcsZt9ZNKD.
 - Thomas Heap et al. Sparse autoencoders can interpret randomly initialized transformers. *arXiv* preprint arXiv:2501.17727, 2025. URL https://arxiv.org/abs/2501.17727.
 - Adam Karvonen, Can Rager, Johnny Lin, Curt Tigges, Joseph Bloom, David Chanin, Yeu-Tong Lau, Eoin Farrell, Callum McDougall, Kola Ayonrinde, Matthew Wearden, Arthur Conmy, Samuel Marks, and Neel Nanda. Saebench: A comprehensive benchmark for sparse autoencoders in language model interpretability, 2025. URL https://arxiv.org/abs/2503.09532.
 - Dmitrii Kharlapenko, Stepan Shabalin, Arthur Conmy, and Neel Nanda. Scaling sparse feature circuits for studying in-context learning. In *Sparsity in LLMs (SLLM): Deep Dive into Mixture of Experts, Quantization, Hardware, and Inference*, 2025. URL https://openreview.net/forum?id=sdLwJTtKpM.
 - Jack Lindsey, Adly Templeton, Jonathan Marcus, Thomas Conerly, Joshua Batson, and Christopher Olah. Sparse crosscoders for cross-layer features and model diffing, 2024. URL https://transformer-circuits.pub/2024/crosscoders/index.html.
 - Scott C. Lowe et al. Logical activation functions: Logit-space equivalents of probabilistic boolean operators. arXiv preprint arXiv:2110.11940, 2021. URL https://arxiv.org/abs/2110.11940.
 - Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. Sparse feature circuits: Discovering and editing interpretable causal graphs in language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=I4e82CIDxv.
 - Anish Mudide, Joshua Engels, Eric J Michaud, Max Tegmark, and Christian Schroeder de Witt. Efficient dictionary learning with switch sparse autoencoders. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=k2ZVAzVeMP.
 - Gonçalo Paulo, Alex Mallen, Caden Juang, and Nora Belrose. Automatically interpreting millions of features in large language models, 2024. URL https://arxiv.org/abs/2410.13928.
 - Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the finest text data at scale, 2024. URL https://arxiv.org/abs/2406.17557.
 - Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Tom Lieberum, Vikrant Varma, János Kramár, Rohin Shah, and Neel Nanda. Improving sparse decomposition of language model activations with gated sparse autoencoders. In *NeurIPS*, 2024a. URL http://papers.nips.cc/paper_files/paper/2024/hash/01772a8b0420baec00c4d59fe2fbace6-Abstract-Conference.html.
 - Senthooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse autoencoders, 2024b. URL https://arxiv.org/abs/2407.14435.
 - Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc V. Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations (ICLR)*, 2017. URL https://arxiv.org/abs/1701.06538.

595

596

597

598

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626 627

629

630

631

632

633 634

635

636

637

638

639

640

641

642

644

645

646

647

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Cogan, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. Gemma 2: Improving open language models at a practical size, 2024. URL https://arxiv.org/abs/2408.00118.

Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024. URL https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html.

Peihao Wang, Rameswar Panda, Lucas Torroba Hennigen, Philip Greengard, Leonid Karlinsky, Rogerio Feris, David Daniel Cox, Zhangyang Wang, and Yoon Kim. Learning to grow pretrained models for efficient transformer training. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=cDYRS5iZ16f.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin

Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.

A ADDITIONAL DETAILS AND RESULTS

A.1 EXPERIMENTAL SETUP

Training details. All SAEs are optimized using AdamW with an initial learning rate of 8×10^{-4} , a cosine learning-rate schedule with a minimum LR of 1×10^{-6} , and a linear warmup for the first 10% of total training steps, auxiliary loss penalty equal to 0.03125. We use a global batch size of 8,192. We sweep over dictionary (latent) sizes of 2^{15} , 2^{16} , and 2^{17} features. For our KronSAE variant, we further sweep the number of heads h and the per-head dimensions m and n such that $h \cdot m \cdot n$ equals the desired dictionary size. Regularization weights and auxiliary loss coefficients are kept constant throughout the runs to isolate the impact of architectural choices.

For all experiments, we spent about 330 GPU days on NVIDIA H100 80GB GPUs, including preliminary research.

SAE. For all experiments on Qwen-2.5, we train each SAE on activations from layer 14. Also for Pythia-1.4B we use layer 14 and for Gemma-2-2B we take activations from layer 12. For most of our experiments, we use sparsity level of $\ell_0 = 50$ non-zero activations per token.

Initialization. As observed by Gao et al. (2025), initializing the decoder as the transpose of the encoder $(W_{\text{dec}} = W_{\text{enc}}^{\top})$ provides a strong metric improvement. We adopt this strategy within KronSAE by partitioning W_{enc} into h head-wise blocks of shapes $m \times d$ and $n \times d$, denoted $\{P_i, Q_i\}_{i=1}^h$. For each head k, we define its decoded rows via a simple additive composition:

$$C_k[i,j] = P_{k,i} + Q_{k,j}, \quad i = 1, \dots, m, \ j = 1, \dots, n.$$

Finally, flattening the matrices $\{C_k\}$ yields full decoder weight matrix $W_{\text{dec}} \in \mathbb{R}^{F \times d}$.

A.2 FLOPS CALCULATION AND EFFICIENCY

For TopK SAE and KronSAE we compute FLOPs in the following way:

$$FLOPS_{TopK}(d, F, k) = dF + kd,$$

$$FLOPS_{KronSAE}(d, m, n, h, k) = dh(m+n) + mnh + kd \approx dh(m+n) + kd.$$
(6)

We calculate FLOPs for most effective variant of TopK where we perform vector matrix multiplication only for nonzero activations, while encoder still requires dense matrix multiplication.

We have also measured the wallclock time for forward and backward to examine the scaling. Figure 8 reports scaling for different hidden dimension sizes.

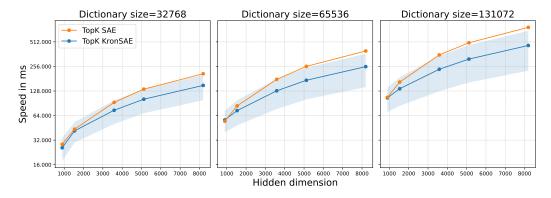


Figure 8: Speed comparision of TopK SAE with KronSAE across different hidden dimensionss. We can see that KronSAE have better scaling properties than SAE with default encoder architecture.

A.3 PYTHIA SUITE

For Pythia-1.4B we train all SAEs on the 12th transformer layer with a budget of 125M tokens. As reported in Table 3, KronSAE achieves performance comparable to TopK SAE with increased number of heads.

Dictionary	SAE	Mean	Max
32k	TopK KronSAE	0.783	0.793 0.793
65k	TopK KronSAE	0.795	0.802 0.805
131k	TopK KronSAE	0.800	0.801 0.810

Table 3: Performance of Pythia-1.4B at 125M budget. At larger dictionary size and fixed training budget KronSAE outperforms TopK SAE.

We conducted additional experiments with smaller Pythia models and trained KronSAE at the middle layers with 512 heads. Table 4 reports results for 125M budget on 65k and 262k dictionary sizes.

Model	70M (d	d=512)	160M (d=768)		410M (d=1024)	
Dictionary	65k	256k	65k	256k	65k	256k
m=1	0.893	0.892	0.856	0.855	0.834	0.835
m=2	0.896	0.897	0.859	0.859	0.832	0.841
m = 4	0.894	0.899	0.857	0.859	0.828	0.839
m = 8	0.896	0.899	0.856	0.862	0.827	0.846
TopK	0.905	0.903	0.870	0.867	0.843	0.847

Table 4: Performance of SAEs on 70M, 160M and 410M Pythias with varying hidden dimensionality.

A.4 COMPARISON WITH JUMPRELU

We provide experiments to compare KronSAE with an alternative activation mechanism, JumpReLU, and report explained variance under three sparsity levels in Table 5.

Model Variant	$\ell_0 = 32$	$\ell_0 = 50$	$\ell_0 = 64$
TopK	0.809	0.837	0.852
JumpReLU	0.813	0.838	0.844
KronSAE (TopK)	0.814	0.840	0.853
KronSAE (JumpReLU)	0.790	0.817	0.828

Table 5: Performance of SAEs with JumpReLU and TopK activations. Since we have floating sparsity controlled via 10 penalty coefficient, we performed a sweep over various sparsity levels, fitted a parabola to the resulting data as a function of ℓ_0 , and evaluated it on those sparsity levels. In contrast, TopK and KronSAE (with TopK) were trained using fixed, predefined sparsity levels.

Replacement of TopK with JumpReLU within KronSAE leads to a degraded performance relative to both JumpReLU SAE and KronSAE with TopK, also with degraded scaling over ℓ_0 . This suggests that the architectural advantages of KronSAE interact most effectively with TopK's hard-thresholding behaviour and its induced sparsity pattern. Whether an alternative activation function can improve on this remains a topic for future research.

A.5 CHOICE OF m, n, h

We derive the following guidelines to train KronSAE: one should minimize the m and maximize the h to improve the reconstruction performance, and search for the most expressive configuration from feasible ones. Since m has more impact on EV, one should start from m=2 in the search process, since it gives improved computational performance with on-par EV with full TopK training.

A.6 ABSORPTION SCORE ON GEMMA AND PYTHIA

In section 4.2 we have analysed whether KronSAE achieves lower absorption score and have answered affirmatively. We also compare the results for Gemma 2 2B and Pythia 1.4B models and validate the improvements, as reported in the Table 6.

Model	SAE	$\ell_0 = 16$	$\ell_0 = 32$	$\ell_0 = 64$	$\mid \ell_0 = 128$	$\ell_0 = 256$
Gemma	TopK	0.410	0.359	0.217	0.047	0.011
	KronSAE	0.040	0.034	0.021	0.008	0.002
Pythia	TopK	0.445	0.233	0.058	0.006	0.004
	KronSAE	0.244	0.129	0.033	0.007	0.003

Table 6: Absorption score calculated for Gemma 2 2B and Pythia 1.4B models. KronSAE shows lower score due to structured latent space and hierarchy between pre-latents and post-latents.

These results confirm that our compositional architectures improves the absorption score and feature consistency across various models from different families.

B More Results on Synthetic

We further evaluate KronSAE on several variant block-diagonal covariance matrices (Figure 9). In each case, the decoder-weight covariance $C_{\rm dec} = W_{\rm dec} W_{\rm dec}^{\top}$ of KronSAE more faithfully reproduces the ground truth groupings than the TopK SAE. Notably, on the third covariance pattern (where some blocks are very small) TopK's learned correlations nearly vanish, whereas KronSAE still uncovers the correct block structure. For the first covariance matrix, KronSAE yields sharply elevated correlations in the regions corresponding to the true blocks, in line with our design goal of head-wise compositionality.

Table 7 quantifies these observations via the RV coefficient and permutation tests. Even after optimally matching TopK's atoms to a dense AE reference, TopK SAE attains only weak correlation alignment ($RV \approx 0.05-0.08$) with non-significant or marginal p-values. In contrast, KronSAE configurations achieve RV values between 0.11 and 0.35 (all p < 0.001), representing a 3–6x improvement in correlation recovery. These results confirm that our compositional encoder not only accelerates training but also robustly captures the intended hierarchical feature interactions across diverse covariance regimes.

To show that KronSAE learn underlying feature covariance better than standard TopK SAE, we change the default KronSAE initialization (which described in Appendix A) to standard normal initialization so to remove bias in results. While RV coefficient for TopK SAE and matched version stays on same level with randomly initialized version, trained KronSAE reveal structure better, as shown in Table 8. We generated data from a purely diagonal covariance matrix (independent features) using the same Cholesky-based procedure. Although fully arbitrary covariance matrices cannot be handled by our Cholesky pipeline, this diagonal covariance matrix can serve as a baseline case with no correlation.

This pivot highlights how KronSAE variants consistently achieve higher RV coefficients, especially the h=2, m=4, n=32 configuration, compared to the TopK baselines across different covariance patterns.

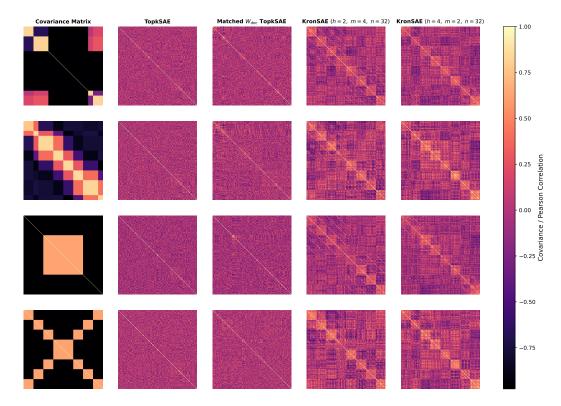


Figure 9: Comparision of how different KronSAE try to reveal hidden structure of defined covariance matrix. The KronSAE models recover the underlying block-structured correlations more faithfully than the TopK baseline, with the finer head/composition split (h=4,m=2) capturing smaller feature groups more accurately.

C MAND AS A LOGICAL OPERATOR

For KronSAE, we replace the original AND_{AIL} (Lowe et al., 2021) with a more restrictive approximation, mAND. Since our objective is to drive each atom toward a distinct, monosemantic feature, we found that tightening the logical conjunction encourages sharper feature separation. Moreover, by using the geometric mean $(\sqrt{p\,q})$ rather than a simple product or minimum, mAND preserves activation magnitudes. A visual comparison of mAND and AND_{AIL} appears in Figure 10.

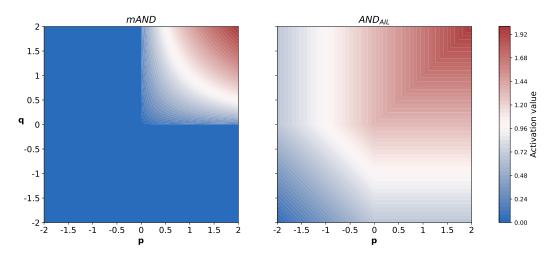


Figure 10: Comparison of the smooth mAND operator against the AND_{AIL} (Lowe et al., 2021).

918	
919	
920	
921	
922	
923	
924	
925	
926	
927	
928	
929	
930	
931	
932	
933	
934	
935	
936	
937	
938	
939	
940	
941	
942	Table '
943	patteri
	o o la i o r

Covariance matrix	SAE	<i>RV</i> p-	-value
	TopK	0.046 0	.1324
1	TopK Matched	0.051 0	.0102
	KronSAE $h = 2, m = 4, n = 32$	0.200 0	.0002
	KronSAE $h = 4, m = 2, n = 32$	0.150 0	.0002
	TopK	0.033 0	.2212
2	TopK Matched	0.080 0	.0002
	KronSAE $h = 2, m = 4, n = 32$	0.221 0	.0002
	$ \overline{ \text{KronSAE } h = 4, m = 2, n = 32 } $	0.319 0	.0002
	ТорК	0.035 0	.2490
3	TopK Matched	0.043 0	.0002
	$ \overline{ \text{KronSAE } h = 2, m = 4, n = 32 } $	0.111 0	.0002
	KronSAE $h = 4, m = 2, n = 32$	0.212 0	.0002
	ТорК	0.034 0	.4613
4	TopK Matched	0.043 0	.0002
	$ \overline{ \text{KronSAE } h = 2, m = 4, n = 32 } $	0.346 0	.0002
	$ \overline{ \text{KronSAE } h = 4, m = 2, n = 32 } $	0.334 0	.0002

Table 7: RV coefficient and permutation p-values for correlation recovery on four synthetic covariance patterns. KronSAE outperforms both the standard and atom-matched TopK SAE by a large margin, achieving statistically significant alignment ($p < 10^{-3}$) across all cases.

Setup / # of covariance matrix	1	2	3	4	5	Mean	Diagonal
ТорК	0.045	0.037	0.040	0.040	0.079	0.048	0.015
TopK Matched	0.059	0.048	0.031	0.061	0.109	0.060	0.015
KronSAE ($h=2, m=4, n=32$)	0.111	0.092	0.069	0.113	0.151	0.107	0.020
KronSAE Init ($h=2,m=4,n=32$)	0.040	0.033	0.029	0.025	0.060	0.037	0.013
KronSAE ($h=4, m=2, n=32$)	0.063	0.062	0.027	0.102	0.127	0.076	0.018
KronSAE Init (h=4,m=2,n=32)	0.042	0.036	0.054	0.047	0.082	0.052	0.027

Table 8: RV coefficient for various experiments, mean value across experiments, and RV value for diagonal matrix (no correlation). Higher values for trained KronSAE variants indicate lower impact of architectural bias, and low values for diagonal matrix and experiment 3 indicate that KronSAE learns correlation when it is actually present in the data.

D FEATURE ANALYSIS METHODOLOGY

We analyze learned features using an established pipeline Bills et al. (2023); Paulo et al. (2024) consisting of three stages: (1) statistical property collection, (2) automatic activation pattern interpretation, and (3) interpretation evaluation. The following subsections detail our implementation.

D.1 DATA COLLECTION

Our collection process uses a fixed-size buffer B=384 per feature, continuing until processing a predetermined maximum token count T_{max} . The procedure operates as follows:

Initial processing batches generate large activation packs of 1M examples, where each example comprises 256-token text segments. When encountering feature activations, we add them to the buffer, applying random downsampling to maintain size B when exceeding capacity. This approach

enables processing arbitrary token volumes while handling rare features that may require extensive sampling.

During collection, we compute online statistics including activation minimums, maximums, means, and frequencies. Post-processing yields two key metrics: token entropy and multitoken ratio. The token entropy is calculated as:

token entropy =
$$-\sum_{i=0}^{n} p_i \cdot \log(p_i)$$
, $p_i = \frac{\text{activations of token i}}{\text{total amount of activations}}$, (7)

where n represents unique activated tokens. The multitoken ratio is:

multitoken ratio =
$$\frac{1}{b} \sum_{i=0}^{b} \frac{\text{number of activations in sequence } i}{\text{total tokens in sequence } i}$$
, (8)

with b < B denoting collected context examples per feature.

We then segment examples using a 31-token context window (15 tokens before/after each activation), potentially creating overlapping but non-duplicated examples. Features with high multitoken ratio may have number of examples significantly exceeding B.

A separate negative examples buffer captures non-activating contexts. Future enhancements could employ predictive modeling (e.g., using frequent active tokens) to strategically populate this buffer with expected-but-inactive contexts, potentially improving interpretation quality.

D.2 FEATURE INTERPRETATIONS

For each feature, we generate interpretations by sampling 16 random activation examples above the median activation quantile and presenting them to Qwen3 14B (Yang et al., 2025) (AWQ-quantized with reasoning enabled). The model produces concise descriptions of the activation patterns. Empirical observations suggest reasoning mode improves interpretation quality, though we lack quantitative measurements. This aligns with findings in (Paulo et al., 2024), which compared standard one-sentence responses with Chain-of-Thought outputs, making model reasoning an interesting direction for future research.

The interpretation process uses the system prompt presented in a Figure 11. User prompts include all special characters verbatim, as some features activate specifically on these characters. A representative (slightly abbreviated) user prompt example is presented on Figure 12.

D.3 EVALUATION PIPELINE

We evaluate interpretations using balanced sets of up to 64 positive (activation quantile > 0.5) and 64 negative examples, employing the same model without reasoning to reduce computational costs. When insufficient examples exist, we maintain class balance by equalizing positive and negative counts. The evaluation uses modified system prompts from (Paulo et al., 2024), with added emphasis on returning Python lists matching the input example count exactly. We discard entire batches if responses are unparseable or contain fewer labels than the number of provided examples.

We calculate two scores.

Detection Score: After shuffling positive/negative examples, we present up to 8 unformatted text examples per batch to the model. The model predicts activations (1/0) for each example, generating up to 128 true/predicted label pairs. The score calculates as:

$$score = \frac{1}{2} \left(\frac{correctly predicted positives}{total positives} + \frac{correctly predicted negatives}{total negatives} \right). \tag{9}$$

Fuzzing Score: We «highlight» activated tokens on sampled examples, from which 50% are correctly labeled positive examples, 25% are mislabeled positive examples, and 25% are randomly labeled negative examples. We present batches of up to 8 examples and the model identifies correct/incorrect labeling, with scoring following Equation 9.

[EXPLANATION]:

105210531054

1055 1056 1057

1076

1077

1078

1079

1026 1027 You are a meticulous AI researcher conducting an important 1028 investigation into patterns found in language. Your task is 1029 to analyze text and provide an explanation that thoroughly 1030 encapsulates possible patterns found in it. 1031 1032 Guidelines: 1033 1034 You will be given a list of text examples on which special 1035 words are selected and between delimiters like «this». If a sequence of consecutive tokens all are important, the entire 1036 sequence of tokens will be contained between delimiters «just 1037 like this». How important each token is for the behavior is 1038 listed after each example in parentheses. 1039 1040 - Your explanation should be a concise STANDALONE PHRASE that 1041 describes observed patterns. 1042 - Focus on the essence of what patterns, concepts and 1043 contexts are present in the examples. 1044 - Do NOT mention the texts, examples, activations or the 1045 feature itself in your explanation. 1046 - Do NOT write "these texts", "feature detects", "the patterns suggest", "activates" or something like that. 1047 - Do not write what the feature does, e.g. instead of 1048 "detects heart diseases in medical reports" write "heart 1049 diseases in medical reports". 1050 - Write explanation in the last line exactly after the 1051

Figure 11: System prompt for feature interpretations.

```
1058
        Examples of activations:
1059
                   Leno«,» a San Francisco Democrat«, said in a
1061
        statement.»'
1062
        Activations:
                      ' said (22.74), statement (27.84), in (27.54)'
1063
        Text: ' city spokesman Tyler Gamble« said in an» email.'
1064
        Activations: ' said (2.92), in (12.81), an (14.91)'
1065
1066
        Text: '
                   towpath at Brentford Lock. «Speaking» on BBC
1067
        London 94'
1068
        Activations: 'Speaking (3.48)'
1069
1070
        Text: ' Michelle, a quadriplegic, « told» DrBicuspid.com'
1071
        Activations: ' told (4.05)'
1072
1073
                   CEO Yves Carcelle« said in a statement».'
        Activations: ' said (19.64), in (29.09), statement (29.39)'
1074
1075
```

Figure 12: Example of user prompt passed to LLM. This feature with 16 examples received the interpretation "Structural elements in discourse, including speech attribution, prepositional phrases, and formal contextual markers" with a detection score of 0.84 and fuzzing score of 0.76.

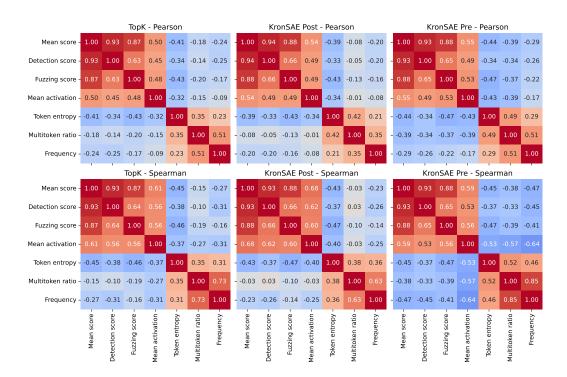


Figure 13: Correlation coefficients (Pearson and Spearman) between properties of TopK and KronSAE latents. Token entropy emerges as a strong predictor of interpretability scores, while higher mean activation and lower frequency also indicate more interpretable features.

E ADDITIONAL FEATURE ANALYSIS RESULTS

Feature property correlations. Our analysis reveals significant correlations between feature properties and interpretability scores (Figure 13). Notably, token entropy and mean activation show substantial correlations with interpretability scores, suggesting their potential as proxies for assessing feature quality without running the full interpretation pipeline. These findings are based on analysis of the first 3072 features from 32k TopK and KronSAE (m=4, n=4) trained on 24M tokens, warranting further validation with larger-scale studies.

Pre-latent to post-latent relationships. We investigate how post-latent properties correlate with various combinations of pre-latent properties, including individual values, means, products, and the mAND operation (product followed by square root). Figure 14 demonstrates that post-latent multitoken ratio, token entropy, and frequency show stronger correlations with pre-latent products or mAND values than with individual pre-latent properties or their means.

Basis geometry. As noted in Section 5.3, latent embeddings primarily exhibit clustering within their originating groups (head, base, extension). With the support of observations reported in Sections ?? and 5.3, we find that models with more heads achieve better reconstruction while producing more diverse basis vectors. This suggests that fine-grained architectures yield more expressive representations, although they may also exhibit undesired challenging behavior like feature splitting (Bricken et al., 2023) or absorption (Chanin et al., 2024).

Figure 15 visualizes this structure through UMAP projections (n_neighbors=15, min_dist=0.05, metric='cosine') of decoder weights from the first 8 heads of 32k SAEs with varying m,n configurations. The plots reveal distinct clustering patterns: for m < n we observe tight base-wise clustering with weaker grouping by extension, and for $m \ge n$ extension-wise clustering is stronger.

This asymmetry suggests that pre-latent capacity requirements directly manifest in the embedding geometry - components with lower polysemanticity (extensions when m < n) exhibit greater geometric

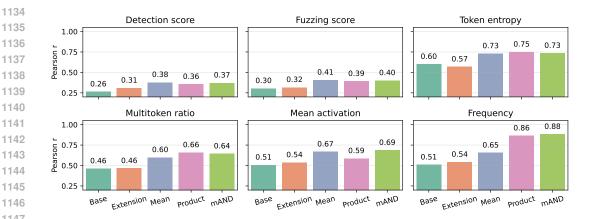


Figure 14: Correlation patterns between properties of post-latents and pre-latents.

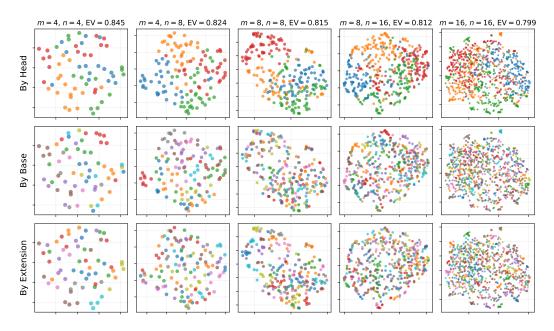


Figure 15: UMAP visualization of post-latent clustering patterns by head, base, and extension group membership. We observe tight clusters by base for m < n and by extension for $m \ge n$.

diversity. We expect symmetric behavior for reciprocal configurations (e.g., m=4,n=8 vs. m=8,n=4), merely swapping the roles of bases and extensions.

KRONSAE IN TERMS OF TENSOR DIAGRAM

The proposed encoder architecture can be visualized as a tensor diagram (Figure 16). Notably, this formulation draws a connection to quantum mechanics, where $|\mathbf{f}\rangle$ represents the (unnormalized) state of two disentangled qubits described by $|p\rangle$ and $|q\rangle$.

If we were to sum the outputs of the encoder's heads instead of concatenating them, $|\mathbf{f}\rangle$ would correspond to a separable quantum state. This scenario can be expressed via the Schmidt decomposition:

$$|\mathbf{f}
angle = \sum_h |oldsymbol{p}_h
angle \otimes_K |oldsymbol{q}_h
angle \, ,$$

where \otimes_K denotes the Kronecker product. However, preliminary experiments revealed that this alternative design results in poorer performance compared to the concatenation-based approach.

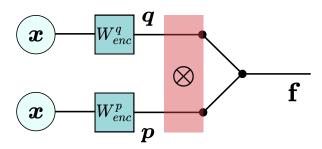


Figure 16: For a single head, the KronSAE encoder architecture separates the input x into two distinct components, p and q, via matrix multiplications with W_{enc}^p and W_{enc}^q accordingly. These components are then combined via the outer product $p \otimes q$, resulting in a matrix representation. To produce the output vector \mathbf{f} , this matrix is flattened into a single vector using a multi-index mapping.

Component	Interpretation	Score			
Extension 2	Scientific instruments, acronyms, and critical numerical values in technical and astronomical contexts				
Base element	s and their compositions with extension 2				
Base 0	<i>Interpretation:</i> Punctuation marks and line breaks serving as structural separators in text.	0.66			
	Composition: Health-related metrics focusing on survival rates, life expectancy, and longevity.	0.88			
Base 1	<i>Interpretation:</i> Numerical values, both in digit form and as spelled-out numbers, often accompanied by punctuation like decimals or commas, in contexts of measurements, statistics, or quantitative expressions.	0.80			
	Composition: Numerical digits and decimal points within quantitative values.	0.86			
Base 2	<i>Interpretation:</i> Nuanced actions and adverbial emphasis in descriptive contexts.	0.71			
	<i>Composition:</i> Astronomical instruments and their components, such as space telescopes and their acronyms, in scientific and observational contexts.	0.90			
Base 3	<i>Interpretation:</i> Forms of the verb "to have" indicating possession, necessity, or occurrence in diverse contexts.	0.91			
	Composition: Antiretroviral therapy components, viral infection terms, and medical treatment terminology.	0.87			

Table 9: Interactions between extension 2 in head 3 and all base elements in that head.

G ANALYSIS OF COMPOSITIONAL STRUCTURE

Here we analyze more examples of interactions in various heads.

Head 3. For this head we have selected all base elements and extension 2, shown in Table 9. Extension element 2 shows moderate interpretability with clear AND-like interactions: with base 1 (semantic inheritance through shared pre-latent semantics) and base 2 (retaining only instrument-related semantics). Notable interactions occur with base 0 (acquiring medical semantics while preserving metric/number aspects) and base 3 (combining instrument semantics with necessity to yield therapy/treatment concepts). The high interpretability scores suggest potential additional encoding mechanisms beyond simple intersection, possibly related to activation magnitude, though dataset or interpretation artifacts cannot be ruled out without further validation.

Component	Interpretation	Score
Extension 2	Hierarchical scopes, geographic references, and spatial dispersal terms	0.78
Base elements	s and their compositions with extension 2	
Base 0	<i>Interpretation:</i> Numerical decimal digits in quantitative expressions and proper nouns.	0.79
	Composition: The state of Illinois in diverse contexts with high significance.	0.95
Base 1	<i>Interpretation:</i> The number three and its various representations, including digits, Roman numerals, and related linguistic forms.	0.84
	Composition: Geographic place names and their linguistic variations in textual contexts.	0.91
Base 2	<i>Interpretation:</i> Ordinal suffixes and temporal markers in historical or chronological contexts.	0.87
	<i>Composition:</i> Terms indicating layers, degrees, or contexts of existence or operation across scientific, organizational, and conceptual domains.	0.82
Base 3	<i>Interpretation:</i> Question formats and topic introductions with specific terms like "What", "is", "of", "the", "Types", "About" in structured text segments.	0.77
	Composition: Spatial spread and occurrence of species or phenomena across environments.	0.87

Table 10: Interactions between extension 2 in head 136 and all base elements in that head.

Component	Interpretation	Score
Extension 1	Geographical mapping terminology and institutional names, phrases involving spatial representation and academic/organizational contexts	
Base element	s and their compositions with extension 1	
Base 0	Interpretation: Proper nouns, abbreviations, and specific named entities. Composition: Geographical or spatial references using the term "map".	0.64 0.93
Base 1	Interpretation: Emphasis on terms indicating feasibility and organizations. Composition: Specific organizations and societal contexts.	0.80 0.89
Base 2	<i>Interpretation:</i> Institutional names and academic organizations, particularly those containing "Institute" or its abbreviations, often paired with prepositions like "of" or "for" to denote specialization or affiliation.	0.89
	Composition: Institutional names containing "Institute" as a core term, often followed by prepositions or additional descriptors.	0.92
Base 3	Interpretation: Closure and termination processes, initiating actions.	0.79
	<i>Composition:</i> Initiating or establishing a state, direction, or foundation through action.	0.85

Table 11: Interactions between extension 1 in head 177 and all base elements in that head.

Head 136. This head exhibits higher interpretability in post-latents than pre-latents. Key observations from the Table 10 include: extension 2 with base 0 narrows semantics to Illinois (likely inheriting geographical subsemantics), while interactions with bases 2-3 demonstrate complexity beyond simple intersection, often introducing additional semantics requiring deeper investigation.

Head 177. Latents presented in Table 11 emonstrates more consistent AND-like behavior than Heads 3 and 136, closely matching the interaction pattern shown in Figure 2.

1297

H KRONSAE SIMPLIFIED IMPLEMENTATION

```
1299 1 class KronSAE (nn.Module):
1300 2
            def __init__(self, config):
                 super().__init__()
1301 3
1302 4
                 self.config = config
                 _t = torch.nn.init.normal_(
1303 6
                          torch.empty(
1304 7
                               self.config.act_size,
1305 8
                               self.config.h * (self.config.m + self.config.n)
                          )
1306 9
1307 10
                      ) / math.sqrt(self.config.dict_size * 2.0)
1308 11
                 self.W_enc = nn.Parameter(_t)
                 self.b_enc = nn.Parameter(
1309 13
                     torch.zeros(self.config.h * (self.config.m + self.config.n))
131014
                 W_dec_v0 = einops.rearrange( # Initialize decoder weights
1311 <sup>15</sup>
                      _t.t().clone(), "(h mn) d -> h mn d",
1312<sup>16</sup>
1312
1313
18
                      h=self.config.h, mn=self.config.m + self.config.n
                 )[:, :self.config.m]
1314<sub>19</sub>
                 W_dec_v1 = einops.rearrange(
131520
                      _t.t().clone(), "(h mn) d -> h mn d",
1316<sup>21</sup>
                      h=self.config.h, mn=self.config.m + self.config.n
1317<sup>22</sup>
                 )[:, self.config.m:]
                 self.W_dec = nn.Parameter(einops.rearrange(
    23
1318 24
                      W_dec_v0[..., None, :] + W_dec_v1[..., None, :, :],
1319<sub>25</sub>
                      "h m n d -> (h m n) d"
1320 26
                 ))
                 self.W_dec.data[:] = (
1321 27
                      self.W_dec.data / self.W_dec.data.norm(dim=-1, keepdim=True)
1322^{28}
1322 <sub>29</sub> 1323 <sub>30</sub>
                 self.b_dec = nn.Parameter(torch.zeros(self.config.act_size))
1324 31
            def encode(self, x: torch.Tensor) -> torch.Tensor:
1325 32
1326<sup>33</sup>
                B, D = x.shape
                 acts = F.relu(
1327 34
                     x @ self.W_enc + self.b_enc
1328 36
                 ).view(B, self.h, self.m + self.n)
1329<sub>37</sub>
                 all scores = torch.sgrt(
1330 38
                      acts[..., :self.config.m, None] * \
1331<sup>39</sup>
                      acts[..., self.config.m:, None, :] + 1e-5
                 ).view(B, -1)
1332<sup>40</sup>
                 scores, indices = all_scores.topk(
1333 42
                      self.config.k, dim=-1, sorted=False
1334<sub>43</sub>
                 )
1335 44
                 acts_topk = torch.zeros(
                     (B, self.config.dict_size)
1336<sup>45</sup>
1337 46
                 ).scatter(-1, indices, scores)
1338 <sub>48</sub>
                 return acts_topk
1339<sub>49</sub>
            def forward(self, x):
                 acts_topk = self.encode(x)
1340 50
                 x_rec = acts_topk @ self.W_dec + self.b_dec
1341 51
1342<sup>52</sup>
                 output = self.get_loss_dict(x, x_rec)
1343 53
1343 54
                 return output
1344<sub>55</sub>
            def get_loss_dict(self, x, x_rec):
1345 56
                 loss = (x_rec - x.pow(2).mean()
                 pt_12 = (x_rec - x).pow(2).sum(-1).squeeze()
1346 57
                 var = (x - x.mean(0)).pow(2).sum(-1).squeeze()
1347 58
1348 <sup>59</sup> <sub>60</sub>
                 ev = (1 - pt_12 / var).mean()
                 return loss, ev
1349
```

I USAGE OF LARGE LANGUAGE MODELS

We have used LLMs as the main tool for conducting the interpretability experiments, as described in section D, and as the instrument for language polishing and word choice.