

---

# Towards Interpretable Scientific Foundation Models: Sparse Autoencoders for Disentangling Dense Embeddings of Scientific Concepts

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The prevalence of foundation models in scientific applications motivates the need  
2 for interpretable representations and search of scientific concepts. In this work, we  
3 present a novel approach using sparse autoencoders (SAEs) to disentangle dense  
4 embeddings from large language models, offering a pathway towards more inter-  
5 pretable scientific foundation models. By training SAEs on embeddings of over  
6 425,000 scientific paper abstracts spanning computer science and astronomy, we  
7 demonstrate their effectiveness in extracting interpretable features while maintain-  
8 ing semantic fidelity. Our method reveals and analyzes SAE features that directly  
9 correspond to scientific concepts, and introduces a novel method for identifying  
10 ‘families’ of related concepts at varying levels of abstraction. To illustrate the  
11 practical utility of our approach, we demonstrate how interpretable features from  
12 SAEs can precisely steer semantic search over scientific literature, allowing for  
13 fine-grained control over query semantics. This work not only bridges the gap  
14 between the semantic richness of dense embeddings and the interpretability needed  
15 for scientific applications, but also offers new directions for improving literature re-  
16 view and scientific discovery. For use by the scientific community, we open-source  
17 our embeddings, trained sparse autoencoders, and interpreted features, along with  
18 a web app for interactive literature search.

## 19 1 Introduction

20 Foundation models have revolutionised natural language processing and are increasingly impacting  
21 scientific domains, enabling the representation of complex scientific concepts in rich semantic spaces  
22 (Devlin et al., 2018; Brown et al., 2020). Dense neural vector embeddings capture nuanced semantic  
23 relationships, enhancing downstream applications such as scientific information retrieval (IR) and  
24 semantic search (Reimers et al., 2019; Gao et al., 2022; Wang et al., 2024). However, the power of  
25 these representations comes at a cost: reduced interpretability and limited user control (Cao et al.,  
26 2023a). In scientific applications, where explainability is critical, these challenges pose a barrier to  
27 embeddings-based tools for literature reviews and scientific discovery.

28 To address these limitations, recent research has explored methods to disentangle and interpret the  
29 information encoded in dense representations (Trifonov et al., 2018). Sparse autoencoders (SAEs)  
30 have emerged as a promising solution for extracting interpretable features from high-dimensional  
31 representations (Ng et al., 2011; Makhzani et al., 2013). By learning to reconstruct inputs as  
32 linear combinations of features in a higher-dimensional sparse basis, SAEs can disentangle complex  
33 representations into individually interpretable components. This approach has shown success in  
34 analysing and steering generative models (Conmy et al., 2024; Lee, 2024; Cunningham et al., 2023b),  
35 but its application to dense text embeddings remains unexplored.

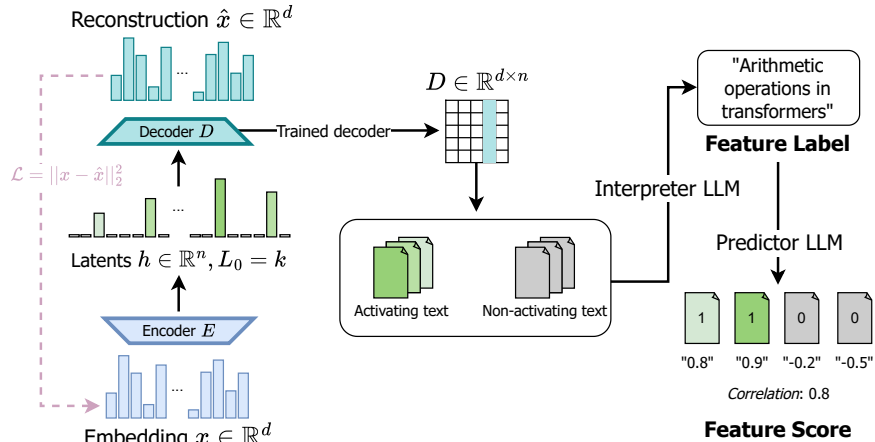


Figure 1: Training and feature labelling process for our sparse autoencoder (SAE). The SAE is trained to minimise reconstruction loss on embeddings from astronomy and computer science paper abstracts. Each feature corresponds to a column in the decoder matrix, representing a direction in embedding space. Feature interpretation involves two steps: (1) An *Interpreter* language model identifies topics present in text that activates each feature but absent in non-activating text. (2) A separate *Predictor* language model assesses feature interpretability by stating its confidence that the feature will activate on unseen text, with confidence correlated with ground truth activations to quantify interpretability.

36 In this work, we present the first application of SAEs to dense text embeddings derived from language  
 37 foundation models, focusing on scientific literature. We demonstrate that this approach offers new  
 38 possibilities for searching, understanding, and manipulating scientific concept spaces. By applying  
 39 our method to embeddings from two diverse scientific domains - computer science and astronomy -  
 40 we showcase its potential for cross-domain applicability in scientific AI. In a direct demonstration  
 41 of their utility for scientific exploration, we show how SAE features can be used to steer scientific  
 42 literature search results, building on previous work applying similar techniques to decoder-only  
 43 transformers and diffusion models for guided generation (Elhage et al., 2022b; Daujotas, 2024). By  
 44 causally manipulating features in the SAE’s hidden representation of an embedding vector, we can  
 45 perform precise adjustments of the semantic meaning of scientific concepts upon reconstruction.

46 Our research makes the following key contributions towards more interpretable scientific foundation  
 47 models. We train SAEs with varying sizes on embeddings from a large corpus of scientific papers  
 48 across two domains, demonstrating their effectiveness in learning interpretable document-level  
 49 features from dense representations of scientific text. We conduct a comprehensive analysis of the  
 50 learned features through the lens of scientific concepts, examining their interpretability, behaviour  
 51 across different model capacities, and semantic properties. To extend this analysis, we introduce the  
 52 concept of SAE “feature families”, clusters of related features that allow for multi-scale analysis and  
 53 manipulation of scientific concepts, and examine how features “split” across levels of abstraction.  
 54 Finally, we demonstrate the practical utility of our approach by applying these interpretable features to  
 55 enhance scientific semantic search, allowing for fine-grained control over query semantics in scientific  
 56 literature exploration. We develop and open-source this as a tool that implements our SAE-enhanced  
 57 semantic search system for scientific literature, as well as open-sourcing the underlying SAEs.

## 58 2 Background and Related work

59 While dense embeddings have dramatically improved performance across various NLP tasks, they  
 60 present significant challenges in terms of interpretability.

### 61 2.1 Embeddings and Representation Learning

62 The evolution of word representations in NLP has progressed from simple one-hot encodings to  
 63 sophisticated dense vector embeddings, culminating in contextual models like BERT (Devlin et al.,

2018) and sentence-level embeddings like Sentence-BERT (Reimers et al., 2019). While these dense embeddings have significantly improved NLP performance, particularly in semantic search and information retrieval (Gao et al., 2021), they present challenges in interpretability and fine-grained control due to their high-dimensional, continuous nature (Liu et al., 2019). This opacity is particularly problematic in applications requiring explainability or precise semantic manipulation. Moreover, dense embeddings face challenges such as the "curse of dense retrieval" (Reimers et al., 2022), where performance degrades rapidly with increasing index size, and difficulties in fine-tuning search results (Cao et al., 2023b; Turian et al., 2010).

## 2.2 Sparse autoencoders

In large language models, the superposition hypothesis suggests that dense neural networks are highly underparameterised, and perform computations involving many more concepts than neurons (Elhage et al., 2022a). Because these semantic concepts, or *features*, are quite sparse, models compensate encoding multiple features within the same set of neurons. However, this also leads to complex, overlapping representations that are difficult to interpret on a single-neuron basis. Similarly, in embedding spaces, features are not represented monosemantically in individual dimensions. Instead, each feature is typically distributed across multiple dimensions, and conversely, each dimension may contribute to the representation of multiple features. This distributed representation allows embedding models to efficiently encode a large number of features in a relatively low-dimensional space, but it also makes the embeddings challenging to interpret directly.

To address this challenge, sparse autoencoders (SAEs) have emerged as a promising solution. SAEs learn to reconstruct inputs using a sparse set of features in a higher-dimensional space, potentially disentangling superposed features (Elhage et al., 2022b; Olshausen et al., 1997). By encouraging this disentanglement, SAEs aim to reveal more interpretable and semantically meaningful representations, demonstrating efficacy in uncovering interpretable features in large language model activations (Donoho, 2006; Gao et al., 2024). In a well-trained SAE, individual features in the hidden dimension align with the underlying sparse, semantically meaningful features.

### 2.2.1 Architecture and training

Sparse autoencoders (SAEs) are neural network models designed to learn compact, interpretable representations of high-dimensional data while enforcing sparsity in the hidden layer activations. The architecture of an SAE consists of an encoder network that maps the input to a hidden representation, and a decoder network that reconstructs the input from this representation.

Let  $\mathbf{x} \in \mathbb{R}^d$  be an input vector, and  $\mathbf{h} \in \mathbb{R}^n$  be the hidden representation, where typically  $n \gg d$ . The encoder and decoder functions are defined as:

$$\text{Encoder : } \mathbf{h} = f_{\theta}(\mathbf{x}) = \sigma(W_e \mathbf{x} + \mathbf{b}_e) \quad (1)$$

$$\text{Decoder : } \hat{\mathbf{x}} = g_{\phi}(\mathbf{h}) = W_d \mathbf{h} + \mathbf{b}_d \quad (2)$$

where  $W_e \in \mathbb{R}^{n \times d}$  and  $W_d \in \mathbb{R}^{d \times n}$  are the encoding and decoding weight matrices,  $\mathbf{b}_e \in \mathbb{R}^k$  and  $\mathbf{b}_d \in \mathbb{R}^d$  are bias vectors, and  $\sigma(\cdot)$  is a non-linear activation function (e.g., ReLU or sigmoid). The parameters  $\theta = \{W_e, \mathbf{b}_e\}$  and  $\phi = \{W_d, \mathbf{b}_d\}$  are learned during training.

The training objective of our SAE combines three main components: a reconstruction loss, a sparsity constraint, and an auxiliary loss. The overall loss function is given by:

$$\mathcal{L}(\theta, \phi) = \frac{1}{d} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \lambda \mathcal{L}_{\text{sparse}}(\mathbf{h}) + \alpha \mathcal{L}_{\text{aux}}(\mathbf{x}, \hat{\mathbf{x}})$$

where  $\lambda > 0$  and  $\alpha > 0$  are hyperparameters controlling the trade-off between reconstruction fidelity, sparsity, and the auxiliary loss.

For the sparsity constraint, we use a  $k$ -sparse constraint: only the  $k$  largest activations in  $\mathbf{h}$  are retained, while the rest are set to zero (Makhzani et al., 2013; Gao et al., 2024). This approach avoids issues such as shrinkage, where L1 regularisation can cause feature activations to be systematically lower than their true values, potentially leading to suboptimal representations *shrinkage*, (Wright et al., 2024; Rajamanoharan et al., 2024). We also use an auxiliary loss, similar to the "ghost grads" technique (Jermyn et al., 2023), to model the reconstruction error using the top  $k_{\text{aux}}$  dead latents, where we typically set  $k_{\text{aux}} = 2k$  (Gao et al., 2024); see Appendix A for details.

## 109 2.2.2 Structure in SAE features

110 State-of-the-art automated interpretability techniques have resulted in the discovery of a large volume  
111 of highly interpretable, monosemantic features in SAEs trained over language models (Cunningham  
112 et al., 2023a; Bricken et al., 2023). With features being the base unit of interpretability for SAEs,  
113 recent work has focused on understanding the geometric structure of features. Bricken et al. (2023)  
114 report *feature splitting* in geometrically close groups of semantically related features, where number  
115 of learned features in the cluster increases with model size. They also report the existence of *universal*  
116 features which re-occur between independent SAEs and which have highly similar activation patterns.  
117 Templeton (2024) find feature splitting also occurs in SAEs trained over production-scale models,  
118 with larger SAEs also exhibiting *novel* features for concepts that are not represented in smaller SAEs.  
119 Makelov et al., 2024 report *over-splitting* of binary features. Engels et al., 2024 find clusters of SAE  
120 features that represent inherently multi-dimensional, non-linear subspaces.

## 121 2.3 Language foundation models in science

122 A number of domain-specific large language models have been developed for question-answering  
123 in specific areas of science, such as medicine or astronomy (Rasmy et al., 2021; Taylor et al., 2022;  
124 Nguyen et al., 2023). Neural vector embeddings from language models have also been leveraged to  
125 enhance scientific question answering and literature search (Kinney et al., 2023; Iyer et al., 2024;  
126 Lála et al., 2023). Recent work has also demonstrated the ability of language models to complete  
127 problem-solving and knowledge synthesis tasks relevant to scientific research (AI4Science et al.,  
128 2023; Romera-Paredes et al., 2024), and even generate novel research hypotheses (Si et al., 2024).  
129 Some research has suggested that some latent scientific knowledge in models exists in structured  
130 representations, and that these representations may be leveraged for scientific discovery (Tshitoyan  
131 et al., 2019; Qu et al., 2024). However, it has also been demonstrated that language models can  
132 exhibit human-like biases or generate false information, limiting their usefulness as scientific tools  
133 (Birhane et al., 2023); here, improved interpretability and steerability could unlock capabilities.

# 134 3 Training SAEs and automated labelling

## 135 3.1 Training and automated interpretability methods

136 We trained top- $k$  Sparse Autoencoders (SAEs) on embeddings of arXiv paper abstracts from astro-  
137 physics (astro-ph, 272,000 papers) and computer science (cs.LG, 153,000 papers) domains, using  
138 OpenAI’s `text-embedding-3-small` model. We experimented with various hyperparameters, fo-  
139 cusing primarily on SAEs with  $k = 16, 32,$  and  $64$  active latents. To interpret the learned features,  
140 we employed an automated two-step process using large language models: an Interpreter to generate  
141 feature labels, and a Predictor to assess interpretation confidence. We evaluated SAEs based on their  
142 reconstruction ability and feature interpretability, using metrics such as normalised mean squared  
143 error and Pearson correlation. Detailed training procedures, hyperparameters, and evaluation metrics  
144 are provided in Appendix A.

145 **SAE Performance:** We observe precise power-law scalings for sparse autoencoder (SAE) perfor-  
146 mance as a function of the number of total latents  $n$ , active latents  $k$ , and compute  $C$  used for  
147 training. The normalised mean squared error (MSE) scales as  $L(n) = cn^{-\alpha}$  for fixed  $k$ , where  $\alpha$   
148 ranges from 0.12 to 0.18, increasing with  $k$ , while  $c$  generally decreases (Figure 2, left panel). For  
149 compute scaling, we calculate the number of training FLOPs  $C$  at each step for each SAE. We find  
150  $L(C) = aC^b$ , where  $a$  generally increases with  $k$  (3.84 for  $k = 16$  to 8.03 for  $k = 64$ ) and  $b$  ranges  
151 from -0.11 to -0.16, becoming more negative as  $k$  increases from 16 to 64 (Figure 2, right panel).  
152 Both relationships show high accuracy with R-squared values above 0.93. Detailed fits are provided  
153 in Appendix A.

154 **Interpretability of SAE features:** The most direct way to evaluate the interpretability of features is  
155 to look at the distribution of automated interpretability scores. Specifically: given a feature label from  
156 our interpreter model, how well can a predictor model predict the feature’s activation on unseen text?  
157 We show in Figure 3 that the Pearson correlation between predictor model confidence of a feature  
158 firing and the ground-truth firing is quite high, with median correlations ranging from 0.65 to 0.71  
159 for cs.LG and 0.85 to 0.98 for astro-ph. We note that Pearson correlation increases as  $k$  and  $n$

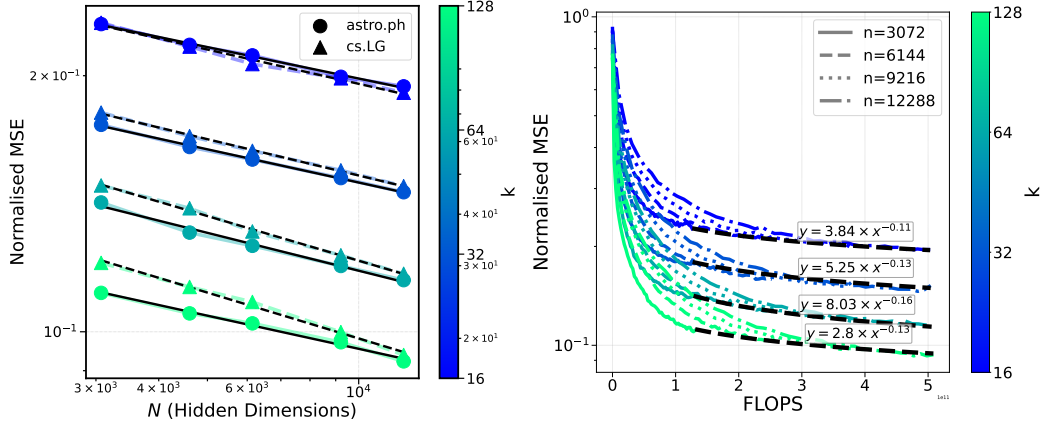


Figure 2: Scaling laws for sparse autoencoder performance. Left: Normalised mean squared error (MSE) as a function of the number of total latents  $n$  for different values of active latents  $k$ . The power-law scaling is evident for each  $k$ . Right: Reconstruction loss as a function of compute (FLOPs) for different  $k$  values, demonstrating the compute-optimal model size scaling.

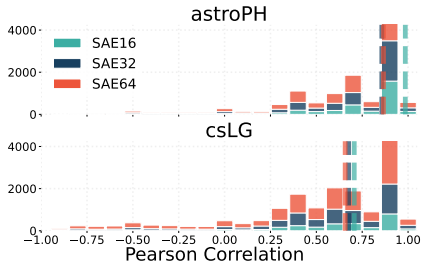


Figure 3: Pearson correlations between the ground-truth and predicted feature activation, using GPT-4o as the *Interpreter* and GPT-4o-mini as the *Predictor*.

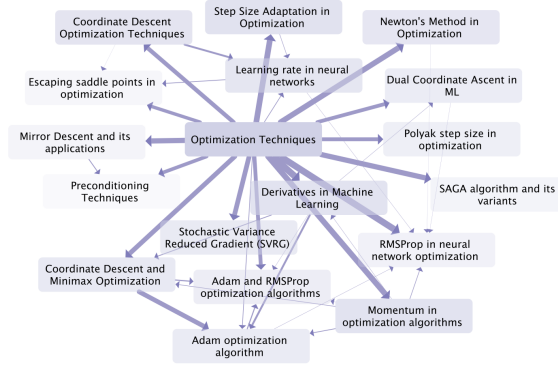


Figure 4: Sample feature family from cs.LG; arrows represent  $C_{ij}^{norm} > 0.1$ , with size  $\propto C_{ij}^{norm}$ .

160 decrease, likely due to models learning coarser-grained features that are easier for the interpreter to  
 161 identify.

## 162 4 Constructing feature families through graph-based clustering

163 We find that our SAEs trained over arXiv paper embeddings recover a wide range of scientifically  
 164 relevant features. These features cover both scientific concepts, from niche to broad and multi-  
 165 disciplinary, and also more abstract semantic artifacts, such as humorous writing or critiques of  
 166 scientific theories. Features and activating examples can be found in Appendix C. In the remainder  
 167 of this work, we focus primarily on features that correspond directly to scientific concepts from  
 168 the literature. Our analysis of feature evolution and grouping provides insights into how scientific  
 169 concepts are represented and related within foundation models, potentially informing the development  
 170 of more interpretable and efficient scientific AI systems.

171 To understand how features evolve across different SAE capacities and to identify meaningful  
 172 groupings of related features, we studied two distinct phenomena: *feature splitting* and *feature*  
 173 *families*. *Feature splitting* – the tendency of features appearing in larger SAEs to “split” the direction  
 174 spanned by a feature from a smaller SAE, and activate on granular sub-topics of the smaller SAE’s  
 175 feature – has been observed in previous work on sparse autoencoders (e.g. Bricken et al., 2023).  
 176 Examples of feature splitting, as well as features recurring across SAEs, can be found in Figures 16  
 177 and 17a/17b.



178 In contrast, *feature families* exist within a single SAE, and exhibit a clear hierarchical structure with  
 179 a dense “parent” feature and several sparser “child” features; we suggest that the “parent” feature  
 180 encompasses a broader, more abstract concept that is shared among the “child” features. An example  
 181 feature family from `cs.LG` can be seen in Figure 4.

## 182 4.1 Feature splitting

183 We investigated how features in smaller SAEs relate to features in larger SAEs through a nearest  
 184 neighbour approach. For each pair of SAEs (i.e. SAE16 and SAE32) with  $n_1$  and  $n_2$  features  
 185 respectively, we calculated an  $n_1 \times n_2$  similarity matrix  $S$  where  $S_{ij} = \mathbf{w}_i^T \mathbf{w}_j / \|\mathbf{w}_i\| \|\mathbf{w}_j\|$ . Here,  
 186  $\mathbf{w}_i$  and  $\mathbf{w}_j$  are decoder weight vectors for features in the smaller and larger SAE, respectively. For  
 187 each feature in the larger SAE, we identified the most similar feature in the smaller SAE, allowing us  
 188 to trace how features potentially “split” or become more refined as model capacity increases.

189 Our results are shown in Figure 9. We find that increasing both number of active latents  $k$  and the  
 190 latent dimension  $n$  reduces the similarity between nearest neighbours in differently sized SAEs. This  
 191 agrees with intuition: larger models with more capacity (higher  $k$  and  $n$ ) can learn more fine-grained  
 192 and specialised features, leading to greater differentiation from features in smaller models.

193 Qualitatively, matching features from small to large SAEs, we find both *recurrent* features and *novel*  
 194 features. Recurrent features appear with extremely high  $S_{ij}$  and activation similarity across one  
 195 or more model pairs, and have highly similar auto-generated interpretations, suggesting semantic  
 196 closeness; these are much more common for lower  $k$  (see Figure 16). In contrast, novel features have  
 197 distinct semantic meaning from their nearest-neighbour match, and activate similarly on some but not  
 198 all documents; novel features thus *split* the semantic space covered by their nearest-neighbour match  
 199 from a smaller SAE. However, some novel features share little semantic or activation overlap with  
 200 their nearest-neighbour feature, as in Fig. 17b, indicating smaller SAEs may not sufficiently cover  
 201 the feature space; see E.1 in the Appendix for more details.

## 202 4.2 Feature families

203 **Feature family identification** To identify feature families, we developed a graph-based approach  
 204 using co-activation patterns across the dataset. We consider only highly interpretable features ( $F1$   
 205  $\geq 0.8$ , Pearson  $\geq 0.8$ ).

206 We first compute co-occurrence matrix  $C$  and activation similarity matrix  $D$ . For all data points  
 207  $k$ ,  $C_{ij} = \sum_k A_{ik} A_{jk}$ ,  $D_{ij} = \sum_k B_{ik} B_{jk}$  where  $A_{ik} = 1$  if feature  $i$  is active on example  $k$  (0  
 208 otherwise), and  $B_{ik} = \mathbf{h}_{k,i}$  if feature  $i$  is active on example  $k$  with hidden vector  $\mathbf{h}_k$  (0 otherwise).  
 209 We normalise the co-occurrence matrix by feature activation frequencies and apply a threshold to  
 210 focus on significant relationships:  $C_{ij}^{norm} = \frac{C_{ij}}{f_i + \epsilon}$  where  $f_i = \sum_k A_{ik}$  is the activation frequency  
 211 of feature  $i$  and  $\epsilon$  is a small constant for numerical stability. We then apply a threshold  $\tau$  to obtain  
 212  $C_{ij}^{thresh}$  (hereafter just  $C$ ). We construct a maximum spanning tree (MST) from  $C$ , capturing the  
 213 strongest relationships between features while avoiding cycles. We convert the MST to a directed  
 214 graph, with edges pointing from higher-density to lower-density features, representing a hierarchy  
 215 from more general to more specific concepts. We identify feature families via depth-first-search  
 216 in this directed graph, starting from root nodes (i.e., no incoming edges) and recursively exploring  
 217 hierarchical sub-families.

218 We iterate this process, removing parent features after each iteration to re-form the MST and  
 219 reveal overlapping, finer-grained feature families. We de-duplicate families with high set overlap  
 220 ( $\frac{|F_1 \cap F_2|}{|F_1 \cup F_2|} > 0.6$ ). In practice, we choose  $\tau = 0.1$  and use  $n = 3$  iterations.

221 **Feature family interpretability** To evaluate the interpretability of feature families and their rele-  
 222 vance to scientific concept understanding, we analysed their collective properties and the effectiveness  
 223 of high-level descriptions in capturing their behaviour in scientific contexts. For each family, we  
 224 generated a “superfeature” description using GPT-4o, based on the individual feature descriptions  
 225 within that family. We then uniformly sampled high-activating examples across all activations of child  
 226 features, and assessed the interpretability of the superfeature using a prediction task, where GPT-4o  
 227 predicted whether test abstracts would activate the superfeature. We compared these predictions to  
 228 ground truth activations to compute Pearson correlation and F1 scores. Additionally, we calculated

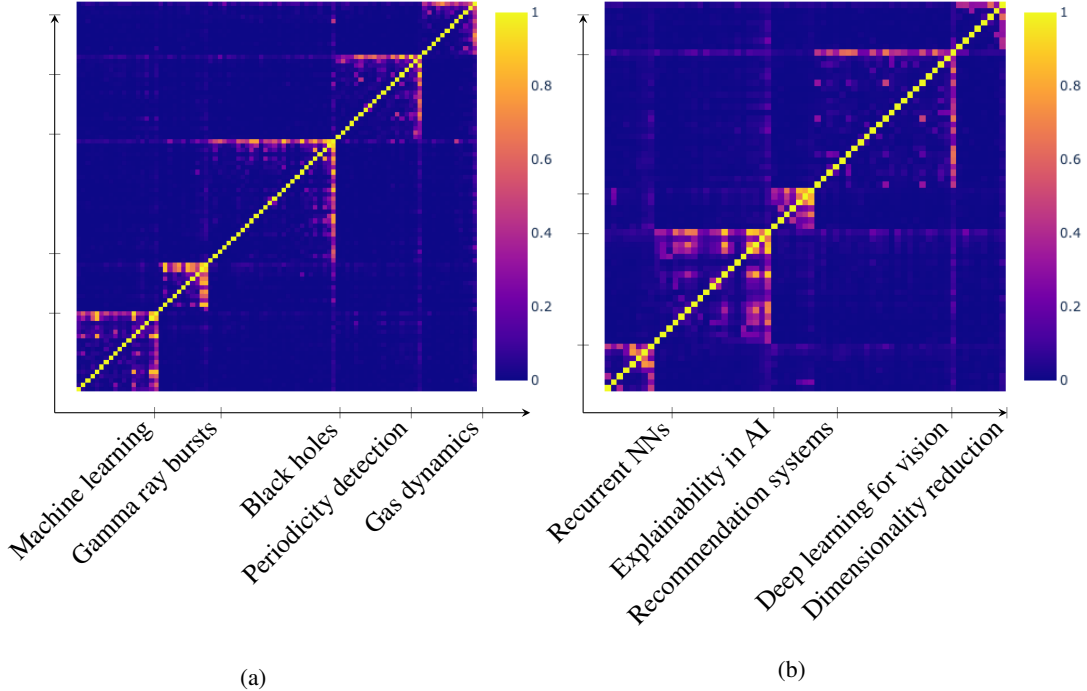


Figure 5: Co-occurrence matrix  $C$  organised by a subset of 5 feature families each. Features in families are ordered by firing density, and the right-most feature is the parent. The un-filled block structure reflects the hierarchical nature of the feature family: all children co-occur with the parent, but few children fire with each other. Visually, this supports our clustering approach.

several metrics to characterise the structure and coherence of the feature families. Table 1 presents the mean values of these metrics across all families for both the `astro-ph` and `cs.LG` datasets.

**Matrix structure** We conjecture that feature families are equivalent to diagonal blocks in some permutation of the co-occurrence matrix  $C$  and activation similarity matrix  $D$ . If feature families are indeed meaningful clusters in the graph, then in  $C$  and  $D$  in-block elements should co-activate much more strongly than off-diagonal elements. We also argue that due to the hierarchical nature of feature families, matrix “blocks” are highly sparse, since child features all co-occur with the parent feature but rarely co-occur with one another. Subsets of the co-occurrence matrix, permuted by feature family, are shown in 5.

Motivated by these structures, we compute the parent-child co-occurrence ratio  $R(p, \mathcal{C})$  for every family with parent feature  $p$  and children  $\mathcal{C}$ ,  $\frac{\text{avg}(\sum_{i \in \mathcal{C}} A_{ip})}{\text{avg}(\sum_{i \in \mathcal{C}} \sum_{j \in \mathcal{C}, j \neq i} A_{ij})}$ . We also permute the co-occurrence and activation similarity matrices by greedily selecting feature families, and compute the in-block to off-diagonal ratios  $C_{\text{diag}}/C_{\text{off}}$  and  $D_{\text{diag}}/D_{\text{off}}$  (excluding the  $i = j$  diagonal), capturing the clustering strength of the block diagonal. Statistics are listed in Table 1.

Dataset	$(k, n)$	Size	F1	Pearson	$\overline{R(p, \mathcal{C})}$	$C_{\text{diag}}/C_{\text{off}}$	$D_{\text{diag}}/D_{\text{off}}$	$f_{\text{inc}}$
<code>astro-ph</code>	(16, 3072)	6	0.86	0.76	10.99	5.13	5.47	0.36
	(32, 6144)	6	0.86	0.73	11.75	4.72	5.87	0.31
	(64, 9216)	7	0.80	0.7	6.87	2.0	3.05	0.24
<code>cs.LG</code>	(16, 3072)	5	0.73	0.6	2.44	8.35	0.89	0.23
	(32, 6144)	5	0.73	0.59	3.5	7.33	1.07	0.30
	(64, 9216)	7	0.80	0.71	1.22	1.78	2.57	0.41

Table 1: Interpretability and structure metrics for feature families from `astro-ph` and `cs.LG`; we report medians unless otherwise noted.  $f_{\text{inc}}$  refers to the fraction of features belonging to a clean (F1  $\geq 0.8$ , Pearson  $\geq 0.8$ ) feature family.

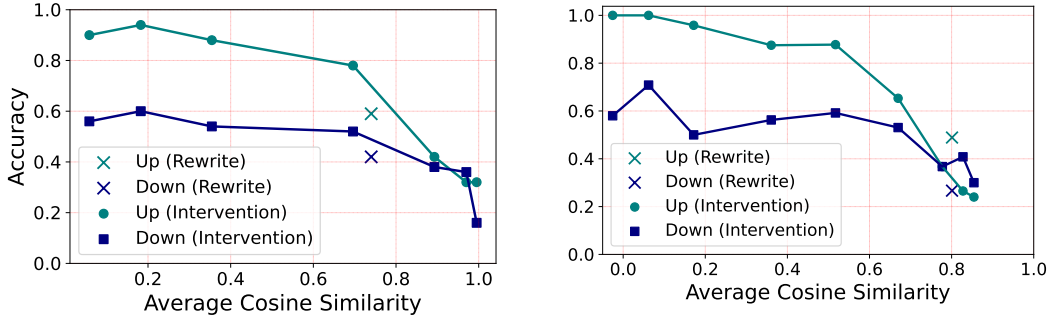


Figure 6: Relationship between intervention accuracy and query fidelity for SAE-based embedding interventions versus traditional query rewriting in scientific literature retrieval for computer science (cs.LG) and astronomy (astro-ph) domains. Intervention accuracy measures the precision of causal query modifications, while query fidelity is quantified by cosine similarity between original and modified query embeddings.

## 243 5 Evaluating effectiveness of search interventions in scientific literature

### 244 5.1 Intervening on scientific embeddings with SAE features

245 As an implementation detail, we note that intervening on a feature by up- or down-weighting its  
 246 hidden representation and then decoding is equivalent to directly adding the scaled feature vector to  
 247 the final embedding. This approach allows for precise manipulation of scientific concepts within the  
 248 embedding space. We explore an alternative process in Appendix G where we iteratively optimise the  
 249 encoded decoded latents to minimise the difference between the desired feature activations and the  
 250 actual activations, potentially offering even finer control over scientific concept representation.

### 251 5.2 Experiments in scientific literature retrieval

252 We incorporate SAE-based embedding interventions into a scientific literature retrieval system for  
 253 computer science (cs.LG) and astronomy (astro-ph), demonstrating cross-domain applicability in  
 254 scientific AI. To assess the effectiveness of SAE feature intervention on semantic search of scientific  
 255 literature, we evaluate the *specificity* and *interpretability* of feature-centric query modifications. We  
 256 select random samples ( $N = 50$  each) of real literature retrieval queries relevant to machine learning  
 257 and astronomy, which are answerable with information in papers from cs.LG and astro-ph.

258 For each scientific query, we return the top  $k = 10$  most relevant papers using embedding cosine  
 259 similarity, forming the original retrieval results  $\mathcal{R}$ . We then select a random feature  $i$  in the top- $k$  from  
 260 the query’s hidden representation  $\mathbf{h}_q$ , and another orthogonal feature  $j$  that has no overlap with the  
 261 top- $k$ ; we limit our selection only to features that are highly interpretable (F1 > 0.9, Pearson > 0.9).  
 262 Given these features, we create a modified query embedding with  $\mathbf{h}'_q, \mathbf{i} = \lambda_-$  and  $\mathbf{h}'_q, \mathbf{j} = \lambda_+$ ,  
 263 letting  $\lambda_- = 0$  and sampling  $\lambda_+ \in [0, 5]$ . This effectively down-weights” and up-weights” the  
 264 importance of specific scientific concepts  $i$  and  $j$ , respectively, in the modified query, which is used  
 265 to generate new retrieval results  $\mathcal{R}'$ .

266 To evaluate the effect of up-weighting and down-weighting query modifications on the final retrieval  
 267 results, we provide both  $\mathcal{R}$  and  $\mathcal{R}'$  to an external LLM instance. The external LLM then compares  
 268  $\mathcal{R}$  and  $\mathcal{R}'$  and determines which scientific concepts, out of a multiple-choice subset of 5 options,  
 269 have been up-weighted or down-weighted; we use this to compute the intervention accuracy, which  
 270 measures the precision and efficacy of causal query interventions in scientific literature search. As a  
 271 baseline, we compare our SAE-based method against traditional query rewriting, by using another  
 272 LLM instance to re-write the original query such that it up-weights  $j$  and down-weights  $i$  entirely  
 273 using natural language. Our results, shown in Figure 6, demonstrate that SAE feature interventions  
 274 consistently outperform traditional query rewriting across various levels of query fidelity in scientific  
 275 literature search.

276 We also experiment with intervening on feature families, sampling highly interpretable families  
 277 containing features in the query top- $k$ . This allows us to manipulate scientific concepts at different



278 levels of abstraction. We uniformly adjust weights for all features in the family, including the parent,  
279 using the auto-generated family interpretation as the multiple-choice option. Results show that  
280 feature family interventions achieve accuracy comparable to individual features, but only down-  
281 weighting interventions outperform query re-writing. This may be because feature families can  
282 comprehensively down-weight related scientific concepts, while up-weighting a general concept  
283 doesn't necessarily require activating all granular child features. Notably, lower cosine similarity  
284 isn't inherently undesirable, as changing the query will naturally reduce similarity.

## 285 6 Discussion

286 In this work, we have presented a novel approach towards more interpretable scientific foundation  
287 models and literature search, by applying sparse autoencoders (SAEs) to dense text embeddings to  
288 derived from large language models. We have demonstrated the usefulness of SAEs in disentangling  
289 embeddings of scientific paper abstracts into semantically relevant document-level concepts, an  
290 important step towards more transparent and controllable AI systems for scientific applications. We  
291 introduced the concept of "feature families" in SAEs, which allow for multi-scale semantic analysis  
292 and manipulation of scientific concepts. Furthermore, we showcased the practical utility of our  
293 approach by applying these interpretable features to enable fine-grained control over query semantics  
294 in scientific literature search, aligning with recent work on controlled text generation (Lee, 2024).

295 Our approach offers a novel solution to the growing challenge of scientific literature exploration. With  
296 the exponential growth in papers, traditional search methods are becoming increasingly ineffective  
297 (Tsang et al., 2016). Our SAE-based approach, for which we provide an open-source interface,  
298 provides a new way to navigate and find pertinent scientific papers, especially in interdisciplinary  
299 fields where relevant work may not be easily discoverable through conventional keyword searches,  
300 citation networks, or vector search (Sharma et al., 2022; Thomsett-Scott et al., 2016).

301 Foundation models, including large language models, are increasingly useful in scientific discovery  
302 (Si et al., 2024; Tshitoyan et al., 2019; AI4Science et al., 2023). By providing concept-level  
303 interpretability, our work also allows for probing the evolution of scientific fields over time, as  
304 captured through state-of-the-art language models and scientific literature corpora. Existing efforts to  
305 map the landscape of scientific research and understand domain and conceptual shifts have relied  
306 primarily on citation networks and keyword analysis (Boyack et al., 2005; Uzzi et al., 2013). However,  
307 SAE features more directly probe semantic meaning and are less sensitive to paper-level or keyword-  
308 level variations, potentially enabling more robust literature searches and meta-analyses. Statistics  
309 of SAE features representing scientific concepts—such as clustering patterns, co-occurrences, and  
310 temporal trends— could gain novel insights into how scientific domains have changed and interacted.

311 To more thoroughly evaluate our approach, we would like to collect human user evaluations of  
312 our SAE-based literature search and compare SAE interventions to other user-facing techniques,  
313 e.g. prompt rewriting. We'd like to evaluate our reconstructed embeddings against the original  
314 embeddings using a standard semantic embedding benchmark such as MTEB (Muennighoff et al.,  
315 2022). We'd also like to be able to conduct an evaluation of SAE features against some proxy of  
316 ground-truth features, much like Makelov et al. (2024) propose. For instance, the Unified Astronomy  
317 Thesaurus (Frey et al., 2018) could provide a basis for evaluating individual feature overlap with  
318 astronomy concepts, and even family features as groupings of these individual concepts.

319 **Limitations:** Our work focused on relatively small datasets from specific scientific domains. Al-  
320 though this specificity allowed us to demonstrate the effectiveness of our approach in targeted areas,  
321 future work should investigate how well these methods generalise to larger, more diverse datasets.  
322 Additionally, our automated interpretability process, while effective, does not utilise the full spectrum  
323 of activations, potentially missing nuanced patterns in feature behaviour.

324 The computational requirements for training SAEs on large embedding datasets also present scalability  
325 challenges that need to be addressed for wider adoption of this approach. Our SAEs are quite small in  
326 comparison to more general language model SAEs. This proved adequate given that we only require  
327 a single embedding vector per example (rather than one per token position) and the narrow domains  
328 we trained on, but SAEs for general text embeddings would need to be scaled up by at least 2-3 the  
329 total number of latents. Further, while we've demonstrated the utility of our approach for literature  
330 search, further work is needed to integrate these interpretable representations into the real-world  
331 workflows of human scientists, from hypothesis generation to experimental design and analysis.

## 332 References

- 333 AI4Science, Microsoft Research and Microsoft Quantum (2023). “The Impact of Large Language  
334 Models on Scientific Discovery: a Preliminary Study using GPT-4”. In: *ArXiv abs/2311.07361*.  
335 URL: <https://api.semanticscholar.org/CorpusID:265150648>.
- 336 Birhane, Abeba, Atoosa Kasirzadeh, David Leslie, and Sandra Wachter (2023). “Science in the age  
337 of large language models”. In: *Nature Reviews Physics* 5, pp. 277–280. DOI: [10.1038/s42254-  
338 023-00581-4](https://doi.org/10.1038/s42254-023-00581-4).
- 339 Boyack, Kevin W, Richard Klavans, and Katy B"orner (2005). “Mapping the backbone of science”.  
340 In: *Scientometrics* 64, pp. 351–374.
- 341 Bricken, Trenton, Catherine Olsson, and Neel Nanda (2023). “Towards Monosemanticity: Decompos-  
342 ing Language Models With Dictionary Learning”. In: *arXiv preprint arXiv:2301.05498*.
- 343 Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
344 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. (2020). “Language  
345 Models are Few-Shot Learners”. In: *Advances in neural information processing systems* 33,  
346 pp. 1877–1901.
- 347 Cao, Wenqiang, Qing Li, Siying Zhang, Rixin Xu, and Youqi Li (2023a). “STEP: Generating Semantic  
348 Text Embeddings with Prompt”. In: *2023 Eleventh International Conference on Advanced Cloud  
349 and Big Data (CBD)*, pp. 180–185. URL: [https://api.semanticscholar.org/CorpusID:  
350 269628678](https://api.semanticscholar.org/CorpusID:269628678).
- 351 Cao, Yichen, Xinyi Wang, Yiran Cao, Renfeng Xu, Zhihan Dong, Qi Fang, Yeyun Gong, Lei Li,  
352 Shuming Shi, Jiafeng Yan, et al. (2023b). “Step-gs: Guiding large language models via step-by-step  
353 prompting”. In: *arXiv preprint arXiv:2305.11725*.
- 354 Conmy, Arthur and Neel Nanda (2024). *Activation Steering with SAEs*. Accessed 16-07-2024. URL:  
355 [https://www.lesswrong.com/posts/C5KAZQib3bzzpeyrg/full-post-progress-  
356 update-1-from-the-gdm-mech-interp-team#Activation\\_Steering\\_with\\_SAEs](https://www.lesswrong.com/posts/C5KAZQib3bzzpeyrg/full-post-progress-update-1-from-the-gdm-mech-interp-team#Activation_Steering_with_SAEs).
- 357 Cunningham, Hoagy, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey (2023a). *Sparse  
358 Autoencoders Find Highly Interpretable Features in Language Models*. arXiv: 2309.08600  
359 [cs.LG]. URL: <https://arxiv.org/abs/2309.08600>.
- 360 – (2023b). “Sparse autoencoders find highly interpretable features in language models”. In: *arXiv  
361 preprint arXiv:2309.08600*.
- 362 Daujotas, Gytis (2024). *Interpreting and Steering Features in Images*. [https://www.lesswrong.  
363 com/posts/Quqekpvx8BGMMcaem/interpreting-and-steering-features-in-images](https://www.lesswrong.com/posts/Quqekpvx8BGMMcaem/interpreting-and-steering-features-in-images).  
364 [Accessed 16-07-2024].
- 365 Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). “BERT: Pre-  
366 training of Deep Bidirectional Transformers for Language Understanding”. In: *arXiv preprint  
367 arXiv:1810.04805*.
- 368 Donoho, David L (2006). “Compressed sensing”. In: *IEEE Transactions on Information Theory* 52.4,  
369 pp. 1289–1306.
- 370 Elhage, Nelson, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,  
371 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish,  
372 Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah (2022a). *Toy Models of  
373 Superposition*. arXiv: 2209.10652 [cs.LG]. URL: <https://arxiv.org/abs/2209.10652>.
- 374 Elhage, Nelson, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Johnston, Ben Mann,  
375 Amanda Askell, Danny Hernandez, Dawn Drain, Zac Hatfield-Dodds, et al. (2022b). “Softmax  
376 Linear Units”. In.
- 377 Engels, Joshua, Isaac Liao, Eric J. Michaud, Wes Gurnee, and Max Tegmark (2024). *Not All Language  
378 Model Features Are Linear*. arXiv: 2405.14860 [cs.LG]. URL: [https://arxiv.org/abs/  
379 2405.14860](https://arxiv.org/abs/2405.14860).
- 380 Frey, Katie and Alberto Accomazzi (2018). “The Unified Astronomy Thesaurus: Semantic metadata  
381 for astronomy and astrophysics”. In: *The Astrophysical Journal Supplement Series* 236.1, p. 24.
- 382 Gao, Leo, John Thickstun, Anirudh Madaan, Zach Scherlis, Arush Guha, Sumanth Dathathri, Jared  
383 Kaplan, Azalia Mirhoseini, and Ilya Sutskever (2024). “Scaling Laws for Neurons in GPT Models”.  
384 In: *arXiv preprint arXiv:2401.02325*.
- 385 Gao, Luyu, Xueguang Ma, Jimmy Lin, and Jamie Callan (2022). *Precise Zero-Shot Dense Retrieval  
386 without Relevance Labels*. arXiv: 2212.10496 [cs.IR]. URL: [https://arxiv.org/abs/  
387 2212.10496](https://arxiv.org/abs/2212.10496).
- 388 Gao, Tianyu, Xingcheng Yao, and Danqi Chen (2021). “SimCSE: Simple contrastive learning of  
389 sentence embeddings”. In: *arXiv preprint arXiv:2104.08821*.

390 Iyer, Kartheik G., Mikael Yunus, Charles O’Neill, Christine Ye, Alina Hyk, Kiera McCormick, Ioana  
391 Ciuca, John F. Wu, Alberto Accomazzi, Simone Astarita, Rishabh Chakrabarty, Jesse Cranney,  
392 Anjalie Field, Tirthankar Ghosal, Michele Ginolfi, Marc Huertas-Company, Maja Jablonska,  
393 Sandor Kruk, Hailing Liu, Gabriel Marchidan, Rohit Mistry, J. P. Naiman, J. E. G. Peek, Mugdha  
394 Polimera, Sergio J. Rodriguez, Kevin Schawinski, Sanjib Sharma, Michael J. Smith, Yuan-Sen  
395 Ting, and Mike Walmsley (2024). *pathfinder: A Semantic Framework for Literature Review  
396 and Knowledge Discovery in Astronomy*. arXiv: 2408.01556 [astro-ph.IM]. URL: <https://arxiv.org/abs/2408.01556>.  
397

398 Jermyn, Adam and Adly Templeton (2023). *Ghost Grads: An improvement on resampling*. [Accessed  
399 19-07-2024]. URL: <https://transformer-circuits.pub/2024/jan-update/index.html#dict-learning-resampling>.  
400

401 Kingma, Diederik P and Jimmy Ba (2014). “Adam: A method for stochastic optimization”. In: *arXiv  
402 preprint arXiv:1412.6980*.

403 Kinney, Rodney, Chloe Anastasiades, Russell Authur, Iz Beltagy, Jonathan Bragg, Alexandra Bu-  
404 raczynski, Isabel Cachola, Stefan Candra, Yoganand Chandrasekhar, Arman Cohan, Miles Craw-  
405 ford, Doug Downey, Jason Dunkelberger, Oren Etzioni, Rob Evans, Sergey Feldman, Joseph  
406 Gorney, David Graham, Fangzhou Hu, Regan Huff, Daniel King, Sebastian Kohlmeier, Bailey  
407 Kuehl, Michael Langan, Daniel Lin, Haokun Liu, Kyle Lo, Jaron Lochner, Kelsey MacMillan, Tyler  
408 Murray, Chris Newell, Smita Rao, Shaurya Rohatgi, Paul Sayre, Zejiang Shen, Amanpreet Singh,  
409 Luca Soldaini, Shivashankar Subramanian, Amber Tanaka, Alex D. Wade, Linda Wagner, Lucy Lu  
410 Wang, Chris Wilhelm, Caroline Wu, Jiangjiang Yang, Angele Zamarron, Madeleine Van Zuylen,  
411 and Daniel S. Weld (2023). *The Semantic Scholar Open Data Platform*. arXiv: 2301.10140  
412 [cs.DL]. URL: <https://arxiv.org/abs/2301.10140>.

413 Lála, Jakub, Odhran O’Donoghue, Aleksandar Shtedritski, Sam Cox, Samuel G. Rodrigues, and  
414 Andrew D. White (2023). *PaperQA: Retrieval-Augmented Generative Agent for Scientific Research*.  
415 arXiv: 2312.07559 [cs.CL]. URL: <https://arxiv.org/abs/2312.07559>.

416 Lee, Linus (2024). *Prism: mapping interpretable concepts and features in a latent space of language*.  
417 Accessed 16-07-2024. URL: <https://thesehist.com/posts/prism>.

418 Liu, Nelson F, Matt Gardner, Yonatan Belinkov, Matthew E Peters, and Noah A Smith (2019).  
419 “Linguistic knowledge and transferability of contextual representations”. In: *arXiv preprint  
420 arXiv:1903.08855*.

421 Makelov, Aleksandar, George Lange, and Neel Nanda (2024). “Towards principled evaluations of  
422 sparse autoencoders for interpretability and control”. In: *arXiv preprint arXiv:2405.08366*.

423 Makhzani, Alireza and Brendan Frey (2013). “K-sparse autoencoders”. In: *arXiv preprint  
424 arXiv:1312.5663*.

425 Muennighoff, Niklas, Nouamane Tazi, Loic Magne, and Nils Reimers (2022). “MTEB: Massive text  
426 embedding benchmark”. In: *arXiv preprint arXiv:2210.07316*.

427 Nanda, Neel (2023). *Open Source Replication & Commentary on Anthropic’s Dictionary Learn-  
428 ing Paper*. [Accessed 22-07-2024]. URL: [https://www.alignmentforum.org/posts/  
429 fKuugaxt2XLtKASKk/open-source-replication-and-commentary-on-anthropic-s](https://www.alignmentforum.org/posts/fKuugaxt2XLtKASKk/open-source-replication-and-commentary-on-anthropic-s).

430 Ng, Andrew et al. (2011). “Sparse autoencoder”. In: *CS294A Lecture notes*. Vol. 72. 2011, pp. 1–19.

431 Nguyen, Tuan Dung, Yuan-Sen Ting, Ioana Ciucă, Charlie O’Neill, Ze-Chang Sun, Maja Jabłońska,  
432 Sandor Kruk, Ernest Perkowski, Jack Miller, Jason Li, Josh Peek, Kartheik Iyer, Tomasz Rózański,  
433 Pranav Khetarpal, Sharaf Zaman, David Brodrick, Sergio J. Rodríguez Méndez, Thang Bui, Alyssa  
434 Goodman, Alberto Accomazzi, Jill Naiman, Jesse Cranney, Kevin Schawinski, and UniverseTBD  
435 (2023). *AstroLLaMA: Towards Specialized Foundation Models in Astronomy*. arXiv: 2309.06126  
436 [astro-ph.IM]. URL: <https://arxiv.org/abs/2309.06126>.

437 Olshausen, Bruno A and David J Field (1997). “Sparse coding with an overcomplete basis set: A  
438 strategy employed by V1?” In: *Vision Research* 37.23, pp. 3311–3325.

439 Qu, Jiaying, Yuxuan Richard Xie, Kamil M. Ciesielski, Claire E. Porter, Eric S. Toberer, and Elif  
440 Ertekin (2024). “Leveraging language representation for materials exploration and discovery”. In:  
441 *npj Computational Materials* 10, pp. 1–14. DOI: 10.1038/s41524-024-01231-8.

442 Rajamanoharan, Senthoran, Arthur Conmy, Lewis Smith, Tom Lieberum, Vikrant Varma, János  
443 Kramár, Rohin Shah, and Neel Nanda (2024). “Improving dictionary learning with gated sparse  
444 autoencoders”. In: *arXiv preprint arXiv:2404.16014*.

445 Rasmy, Laila, Yang Xiang, Ziqian Xie, Cui Tao, and Degui Zhi (2021). “Med-BERT: pretrained con-  
446 textualized embeddings on large-scale structured electronic health records for disease prediction”.  
447 In: *npj Digital Medicine* 4.1, pp. 1–13.

448 Reimers, Nils, Lucas Beyer, and Iryna Wang (2022). “The curse of dense low-dimensional information  
449 retrieval for large index sizes”. In: *arXiv preprint arXiv:2112.07899*.

450 Reimers, Nils and Iryna Gurevych (2019). “Sentence-BERT: Sentence Embeddings using Siamese  
451 BERT-Networks”. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural  
452 Language Processing*, pp. 3982–3992.

453 Romera-Paredes, Bernardino, Mohammadamin Barekatain, Alexander Novikov, et al. (2024). “Mathe-  
454 matical discoveries from program search with large language models”. In: *Nature* 625, pp. 468–475.  
455 DOI: [10.1038/s41586-023-06924-6](https://doi.org/10.1038/s41586-023-06924-6).

456 Sharma, Ritu, Sarita Gulati, Amanpreet Kaur, Atasi Sinhababu, and Rupak Chakravarty (2022).  
457 “Research discovery and visualization using ResearchRabbit: A use case of AI in libraries”. In:  
458 *COLLNET Journal of Scientometrics and Information Management* 16.2, pp. 215–237.

459 Si, Chenglei, Diyi Yang, and Tatsunori Hashimoto (2024). *Can LLMs Generate Novel Research  
460 Ideas? A Large-Scale Human Study with 100+ NLP Researchers*. arXiv: [2409.04109](https://arxiv.org/abs/2409.04109) [cs.CL].  
461 URL: <https://arxiv.org/abs/2409.04109>.

462 Taylor, Ross, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia,  
463 Andrew Poulton, Viktor Kerkez, and Robert Stojnic (2022). *Galactica: A Large Language Model  
464 for Science*. arXiv: [2211.09085](https://arxiv.org/abs/2211.09085) [cs.CL]. URL: <https://arxiv.org/abs/2211.09085>.

465 Templeton, Adly (2024). *Scaling monosemanticity: Extracting interpretable features from claude 3  
466 sonnet*. Anthropic.

467 Thomsett-Scott, Beth and Patricia E Reese (2016). “Academic libraries and discovery tools: A survey  
468 of the literature”. In: *Discovery Tools: The Next Generation of Library Research*, pp. 3–23.

469 Trifonov, Valentin, Octavian-Eugen Ganea, Anna Potapenko, and Thomas Hofmann (2018). *Learning  
470 and Evaluating Sparse Interpretable Sentence Embeddings*. arXiv: [1809.08621](https://arxiv.org/abs/1809.08621) [cs.CL]. URL:  
471 <https://arxiv.org/abs/1809.08621>.

472 Tsang, Daniel C and Julia M Gelfand (2016). “The Changing Landscape of Research Library  
473 Collections: Ensuring Realistic Sustainability.” In.

474 Tshitoyan, Vahe, John Dagdelen, Leigh Weston, Alexander Dunn, Ziqin Rong, Olga Kononova,  
475 Kristin A. Persson, Gerbrand Ceder, and Anubhav Jain (2019). “Unsupervised word embeddings  
476 capture latent knowledge from materials science literature”. In: *Nature* 571, pp. 95–98. DOI:  
477 [10.1038/s41586-019-1335-8](https://doi.org/10.1038/s41586-019-1335-8).

478 Turian, Joseph, Lev Ratinov, and Yoshua Bengio (2010). “Word representations: a simple and  
479 general method for semi-supervised learning”. In: *Proceedings of the 48th annual meeting of the  
480 association for computational linguistics*, pp. 384–394.

481 Uzzi, Brian, Satyam Mukherjee, Michael Stringer, and Ben Jones (2013). “Atypical combinations  
482 and scientific impact”. In: *Science* 342.6157, pp. 468–472.

483 Wang, Liang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder,  
484 and Furu Wei (2024). *Text Embeddings by Weakly-Supervised Contrastive Pre-training*. arXiv:  
485 [2212.03533](https://arxiv.org/abs/2212.03533) [cs.CL]. URL: <https://arxiv.org/abs/2212.03533>.

486 Wright, Benjamin and Lee Sharkey (2024). *Addressing Feature Suppression in SAEs*. [https://www.  
487 alignmentforum.org/posts/3JuSjTZyMzaSeTxKk/addressing-feature-suppression-  
488 in-saes](https://www.alignmentforum.org/posts/3JuSjTZyMzaSeTxKk/addressing-feature-suppression-in-saes). [Accessed 16-07-2024].

489	<b>Contents</b>	
490	<b>1 Introduction</b>	<b>1</b>
491	<b>2 Background and Related work</b>	<b>2</b>
492	2.1 Embeddings and Representation Learning . . . . .	2
493	2.2 Sparse autoencoders . . . . .	3
494	2.2.1 Architecture and training . . . . .	3
495	2.2.2 Structure in SAE features . . . . .	4
496	2.3 Language foundation models in science . . . . .	4
497	<b>3 Training SAEs and automated labelling</b>	<b>4</b>
498	3.1 Training and automated interpretability methods . . . . .	4
499	<b>4 Constructing feature families through graph-based clustering</b>	<b>5</b>
500	4.1 Feature splitting . . . . .	6
501	4.2 Feature families . . . . .	6
502	<b>5 Evaluating effectiveness of search interventions in scientific literature</b>	<b>8</b>
503	5.1 Intervening on scientific embeddings with SAE features . . . . .	8
504	5.2 Experiments in scientific literature retrieval . . . . .	8
505	<b>6 Discussion</b>	<b>9</b>
506	<b>A Training details</b>	<b>14</b>
507	A.1 Training setup . . . . .	14
508	A.2 Training and automated interpretability methods . . . . .	14
509	A.3 SAE training metrics . . . . .	15
510	A.4 Scaling laws . . . . .	16
511	A.5 Feature density and similarity . . . . .	16
512	<b>B SAErch.ai</b>	<b>18</b>
513	B.1 Overview . . . . .	18
514	B.2 Feature Visualisation Tab . . . . .	19
515	B.2.1 Individual Features . . . . .	19
516	B.2.2 Feature Families . . . . .	19
517	<b>C Automated interpretability details</b>	<b>20</b>
518	C.1 Examples of features . . . . .	20
519	C.2 Automated interpretability prompts . . . . .	21
520	C.3 Exploring the effectiveness of smaller models . . . . .	23
521	<b>D Cross-domain features</b>	<b>23</b>



522	<b>E Feature family details</b>	<b>25</b>
523	E.1 Feature splitting structures . . . . .	25
524	E.2 Feature family structure . . . . .	26
525	E.3 Feature family interpretability . . . . .	29
526	<b>F Exploring learned decoder weight matrices</b>	<b>29</b>
527	<b>G Iterative encoding optimisation</b>	<b>29</b>

## 528 A Training details

### 529 A.1 Training setup

530 Our sparse autoencoder (SAE) implementation incorporates several recent advancements in the field.  
 531 Following Bricken et al. (2023), we initialise the bias  $b_{pre}$  using the geometric median of a data  
 532 point sample and set encoder directions parallel to decoder directions. Decoder latent directions are  
 533 normalised to unit length at initialisation and after each training step. For our top- $k$  models, based on  
 534 Gao et al. (2024), we set initial encoder magnitudes to match input vector magnitudes, though our  
 535 analyses indicate minimal impact from this choice.

536 We also use an auxiliary loss, similar to the “ghost grads” technique (Jermyn et al., 2023), to model  
 537 the reconstruction error using the top  $k_{aux}$  dead latents, where we typically set  $k_{aux} = 2k$  (Gao et al.,  
 538 2024). Latents are flagged as dead during training if they have not activated for a predetermined  
 539 number of tokens (in our case, one full epoch through the training data). Given the reconstruction  
 540 error of the main model  $\mathbf{e} = \mathbf{x} - \hat{\mathbf{x}}$ , we define the auxiliary loss as  $\mathcal{L}_{aux}(\mathbf{x}, \hat{\mathbf{x}}) = \|\mathbf{e} - \hat{\mathbf{e}}\|_2^2$  where  
 541  $\hat{\mathbf{e}} = W_d \mathbf{z}$  is the reconstruction using the top  $k_{aux}$  dead latents, and  $\mathbf{z}$  is the sparse representation  
 542 using only these dead latents. This additional loss term helps to revive dead features and improve  
 543 the overall representational capacity of the model (Gao et al., 2024). We found that dead latents  
 544 only occurred during training the  $k = 16$  models, and all dead latents had disappeared by the end  
 545 of training. We show how dead latents evolved over training the  $k = 16$  SAEs for the astro-ph  
 546 abstracts in Figure 7.

547 For optimisation, we employ Adam (Kingma et al., 2014) with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ , maintaining  
 548 a constant learning rate. We use gradient clipping. Our training uses batches of 1024 abstracts, with  
 549 performance metrics showing robustness to batch size variations under appropriate hyperparameter  
 550 settings.

551 The primary MSE loss uses a global normalisation factor computed at training initiation, while  
 552 the AuxK loss employs per-batch normalisation to adapt to evolving error distributions. Following  
 553 Bricken et al. (2023), we apply a gradient projection technique to mitigate interactions between the  
 554 Adam optimiser and decoder normalisation.

### 555 A.2 Training and automated interpretability methods

556 **Training:** We train our top- $k$  SAEs on the embeddings of abstracts from papers on arXiv with the  
 557 astro-ph tag (astrophysics, 272,000 papers) and the cs.LG tag (computer science, 153,000 papers).  
 558 The embeddings were generated with OpenAI’s text-embedding-3-small model.<sup>1</sup> We train our  
 559 SAEs on these collections of embeddings separately. We normalised the embeddings to zero mean  
 560 and unit variance before passing them to the SAE as inputs. Our trained SAEs will be made available  
 561 for download.

562 **Hyperparameters:** Notable hyperparameters include the number of active latents  $k$ , the total number  
 563 of latents  $n$ , the number of auxiliary latents  $k_{aux}$ , the learning rate, and the auxiliary loss coefficient  
 564  $\alpha$ . We found learning rate and auxiliary loss coefficient to not have a significant effect on final  
 565 reconstruction loss; we set the former to 1e-4 and the latter to 1/32. We vary  $k$  between 16 and 128,  
 566 and  $n$  between two to nine times the embedding dimension  $d_{input}$ . Whilst we train SAEs with many

<sup>1</sup><https://openai.com/index/new-embedding-models-and-api-updates/>



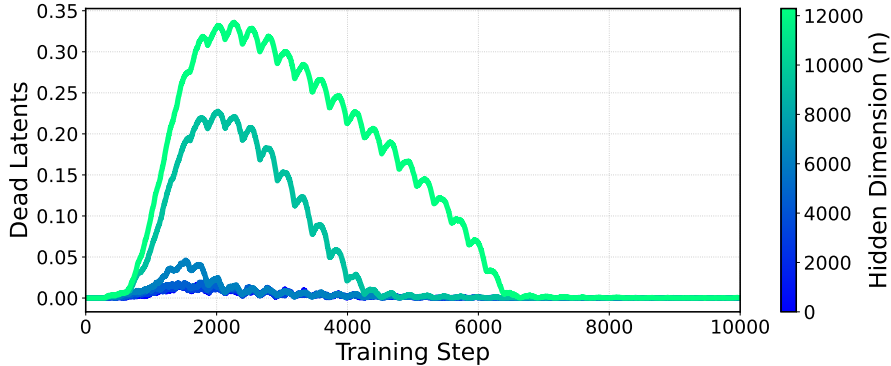


Figure 7: The proportion of dead latents, defined as features that haven’t fired in the last epoch of training, for our  $k = 16$  SAEs on the *astro-ph* abstract embeddings. All dead latents were gone by the end of training. We found that dead latents only occurred in  $k = 16$  autoencoders.

567 different combinations of these hyperparameters, we largely focus on what we hereon refer to as  
 568 SAE16 ( $k = 16$ ,  $n = 2d_{\text{input}} = 3072$ ), SAE32 ( $k = 32$ ,  $n = 4d_{\text{input}} = 6144$ ) and SAE64 ( $k = 64$ ,  
 569  $n = 6d_{\text{input}} = 9216$ ). We train each model for approximately 13.2 thousand steps.

570 **Automated interpretability:** Following the training of a Sparse Autoencoder (SAE), it becomes  
 571 necessary to interpret its features, each corresponding to a column in the learned decoder weight  
 572 matrix. To facilitate feature interpretation and quantify interpretation confidence, we employ two  
 573 Large Language Model (LLM) instances: the *Interpreter* and the *Predictor*. The Interpreter is  
 574 tasked with generating labels for each feature. It is provided with the abstracts that produce the top  
 575 5 activations of the feature across the dataset, along with randomly selected abstracts that do not  
 576 activate the feature. The Interpreter then generates a label for the feature based on this input (for the  
 577 complete prompt, refer to Appendix C). Subsequently, the generated label is passed to the Predictor.  
 578 The Predictor is presented with three randomly sampled abstracts where the feature was activated and  
 579 three where it was not. It is then instructed to predict whether a given abstract should activate the  
 580 feature, expressing its confidence as a score ranging from  $-1$  (absolute certainty of non-activation) to  
 581  $+1$  (absolute certainty of activation).<sup>2</sup> We measure the Pearson correlation between this confidence  
 582 and the true activation (binary;  $+1$  or  $-1$ ). We also measure the F1 score, when framing the confidence  
 583 as a binary classification (active if confidence is above 0, inactive otherwise).

584 **Evaluation metrics:** In order to compare SAEs, we evaluate both their ability to reconstruct the  
 585 embeddings, as well as the interpretability of the learned features. For the former, we examine the  
 586 normalised mean squared error (MSE), where we divide MSE by the error when predicting the mean  
 587 activations. We also report the log density of the activation of features across all papers. We do not  
 588 report dead latents (those not firing on any abstract) as all models contained zero dead latents at the  
 589 end of training. We also report the mean activation of features, when their activation is non-zero. To  
 590 measure interpretability, we use Pearson correlation, as outlined above.

### 591 A.3 SAE training metrics

592 Table 2 shows the final training metrics for all combinations of SAEs trained. We note clear trends in  
 593 normalised MSE, log feature density and activation mean as we vary the number of active latents  $k$   
 594 and the overall number of latents  $n$ .

<sup>2</sup>We use 3 activating and 3 non-activating abstracts for the Predictor, rather than 5, due to LLM costs. We used *gpt-4o* as the Interpreter and *gpt-4o-mini* as the Predictor. Notably, we predict each abstract separately, rather than batching abstracts like Bricken et al. (2023).

Table 2: Metrics for our top- $k$  sparse autoencoders with varying  $k$  and hidden dimensions, across both astronomy and computer science papers. MSE is normalised mean squared error, Log FD is the mean log density of feature activations, and activation mean is the mean activation value across non-zero features. Note that MSE is normalised.

$k$	$n$	astro.ph			cs.LG		
		MSE	Log FD	Act Mean	MSE	Log FD	Act Mean
16	3072	0.2264	-2.7204	0.1264	0.2284	-2.7314	0.1332
	4608	0.2246	-4.7994	0.1350	0.2197	-3.0221	0.1338
	6144	0.2128	-3.1962	0.1266	0.2089	-3.2299	0.1342
	9216	0.1984	-3.4206	0.1264	0.1962	-3.4833	0.1343
	12288	0.1957	-6.2719	0.1274	0.1897	-3.6448	0.1347
32	3072	0.1816	-2.3389	0.0847	0.1831	-2.3008	0.0885
	4608	0.1691	-3.6091	0.0882	0.1697	-2.5152	0.0876
	6144	0.1604	-2.7761	0.0841	0.1641	-2.6687	0.0873
	9216	0.1554	-3.0227	0.0842	0.1540	-2.9031	0.0875
	12288	0.1520	-4.9505	0.0843	0.1457	-3.0577	0.0877
64	3072	0.1420	-1.9538	0.0566	0.1485	-1.8875	0.0584
	4608	0.1331	-2.7782	0.0622	0.1370	-2.0637	0.0570
	6144	0.1262	-2.2828	0.0545	0.1310	-2.1852	0.0558
	9216	0.1182	-2.4682	0.0539	0.1240	-2.3536	0.0545
	12288	0.1152	-3.4787	0.0583	0.1162	-2.4847	0.0548
128	3072	0.1111	-1.8876	0.0483	0.1206	-1.5311	0.0399
	4608	0.1033	-2.1392	0.0457	0.1137	-1.6948	0.0376
	6144	0.1048	-2.2501	0.0438	0.1076	-1.8079	0.0366
	9216	0.0975	-2.5352	0.0409	0.0999	-1.9701	0.0348
	12288	0.0936	-2.7025	0.0399	0.0942	-2.0858	0.0342

#### 595 A.4 Scaling laws

596 For the left panel of Figure 2, which shows the scaling of normalised MSE with the number of total  
 597 latents  $n$ , we observe the following power-law relationships:

$$\begin{aligned}
 k = 16 : L(n) &= 0.61n^{-0.12} \text{ (astro.ph)}; L(n) = 0.67n^{-0.13} \text{ (cs.LG)} \\
 k = 32 : L(n) &= 0.49n^{-0.13} \text{ (astro.ph)}; L(n) = 0.56n^{-0.14} \text{ (cs.LG)} \\
 k = 64 : L(n) &= 0.46n^{-0.15} \text{ (astro.ph)}; L(n) = 0.60n^{-0.17} \text{ (cs.LG)} \\
 k = 128 : L(n) &= 0.31n^{-0.13} \text{ (astro.ph)}; L(n) = 0.51n^{-0.18} \text{ (cs.LG)}
 \end{aligned}$$

598 For the right panel of Figure 2, which shows the scaling of normalised MSE with the amount of  
 599 compute  $C$  (in FLOPs), we observe the following power-law relationships:

$$\begin{aligned}
 k = 16 : L(C) &= 3.84C^{-0.11} \\
 k = 32 : L(C) &= 5.25C^{-0.13} \\
 k = 64 : L(C) &= 8.03C^{-0.16} \\
 k = 128 : L(C) &= 2.80C^{-0.13}
 \end{aligned}$$

600 These equations demonstrate the consistent power-law scaling behaviour of sparse autoencoders  
 601 across different values of  $k$ ,  $n$ , and compute  $C$ .

#### 602 A.5 Feature density and similarity

603 We find an intuitive relationship between  $k$  and  $n$  and the log feature density (essentially, how often a  
 604 given feature fires). As  $k$  increases, we get a sharper peak of log feature density, shifted to the right,  
 605 suggesting features fire in a tighter range as we increase the instantaneous LO of the SAE’s encoder  
 606 (Figure 8).

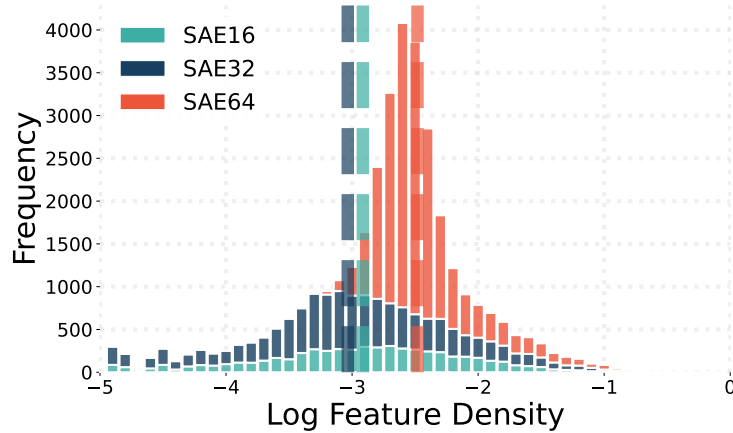
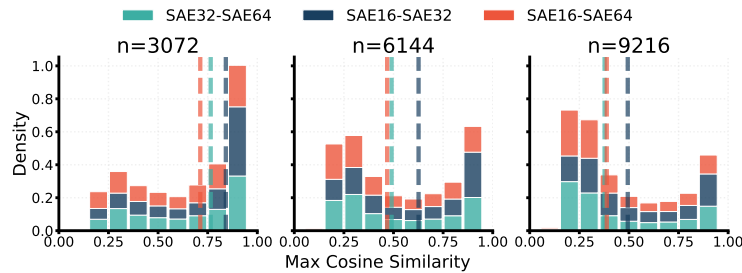
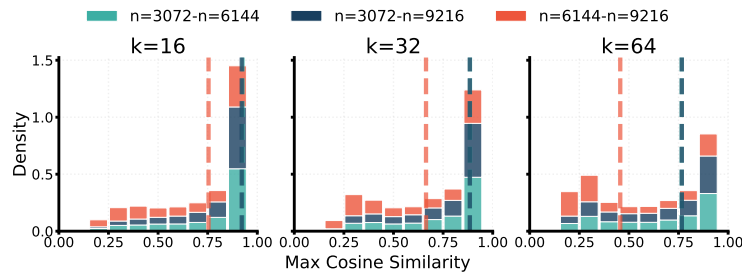


Figure 8: Log feature density for features in our three SAEs as a stacked histogram, showing the distribution of how often features fire across all paper abstracts (*cs.LG* and *astro-ph*). The larger SAE has a higher mean feature density than the smaller SAEs.



(a)  $k$  fixed, varying  $n$ . As  $n$  increases, the features between across SAEs with varying  $k$  become more disparate.



(b)  $n$  fixed, varying  $k$ . Higher values of  $k$  lead to less similarity regardless of  $n$ .

Figure 9: Nearest-neighbour cosine similarity distributions for SAE features. To find features in an SAE with a lower  $k$  that are most similar to those in an SAE with a larger  $k$ , we compute the cosine similarity between each feature in the larger model and each feature in the smaller model. We do this for several values of  $n$ , and combine the distributions for *astro-ph* and *cs.LG*.

607 To compare features across different SAEs trained on the same input data, we analyse the cosine  
 608 similarity between the decoder weight vectors corresponding to each feature. Decoder weights,  
 609 represented by columns in the decoder matrix, directly encode each feature’s contribution to input  
 610 reconstruction. Encoder weights, on the other hand, are optimised to extract feature coefficients  
 611 while minimising interference between non-orthogonal features. This separation is important in the  
 612 context of superposition, where we have more features than input dimensions, precluding perfect  
 613 orthogonality.

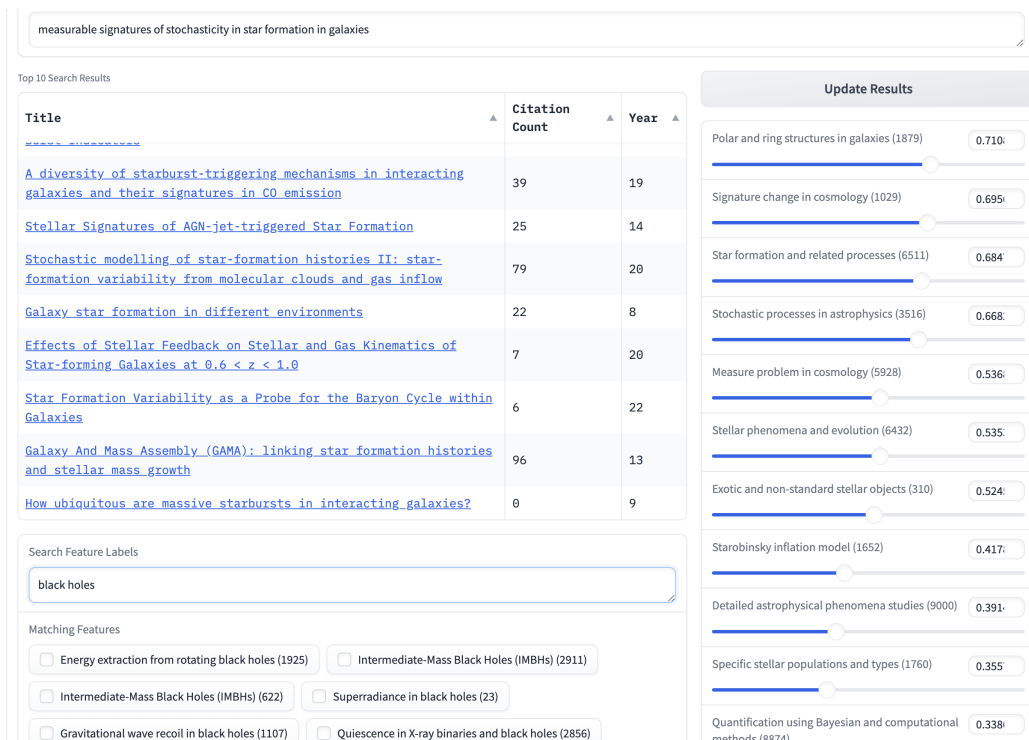


Figure 10: The SAErch tab of our web application, demonstrating a semantic search for “measurable signatures of stochasticity in star formation in galaxies” in the astrophysics domain. The interface displays the top 10 search results ranked by relevance, including title, citation count, and publication year. On the right, sliders represent the top activated SAE features for the query, allowing users to fine-tune the search by adjusting feature weights. On the bottom we have our feature addition interface. Users can search for specific semantic features (e.g., “black holes”) and add them to their query. They can then adjust the strength of these features.

## 614 B SAErch.ai

615 To demonstrate the practical applications of our sparse autoencoder (SAE) approach to semantic  
 616 search and feature interpretation, we developed a web application that allows users to interact with  
 617 the SAE models trained on arXiv paper embeddings. The link will be made public at the end of the  
 618 anonymity period.

### 619 B.1 Overview

620 SAErch.ai is built using the Gradio framework and consists of three main tabs: Home, SAErch, and  
 621 Feature Visualisation. The application allows users to switch between the Computer Science (cs.LG)  
 622 and Astrophysics (astro-ph) datasets.

623 The SAErch tab implements the core functionality of our semantic search system, allowing users to:

- 624 • Input a search query
- 625 • View the top 10 search results based on embedding similarity
- 626 • Interact with the SAE features activated by their query

627 For each query, the system displays sliders corresponding to the top-k SAE features activated by the  
 628 input. Users can adjust these sliders to modify the query embedding, effectively steering the search  
 629 results towards or away from specific semantic concepts; see Figure 10. This directly demonstrates  
 630 the fine-grained control over query semantics discussed in Section 5 of our paper. Users can also  
 631 search for and add specific features not initially activated by their query (Figure ??).

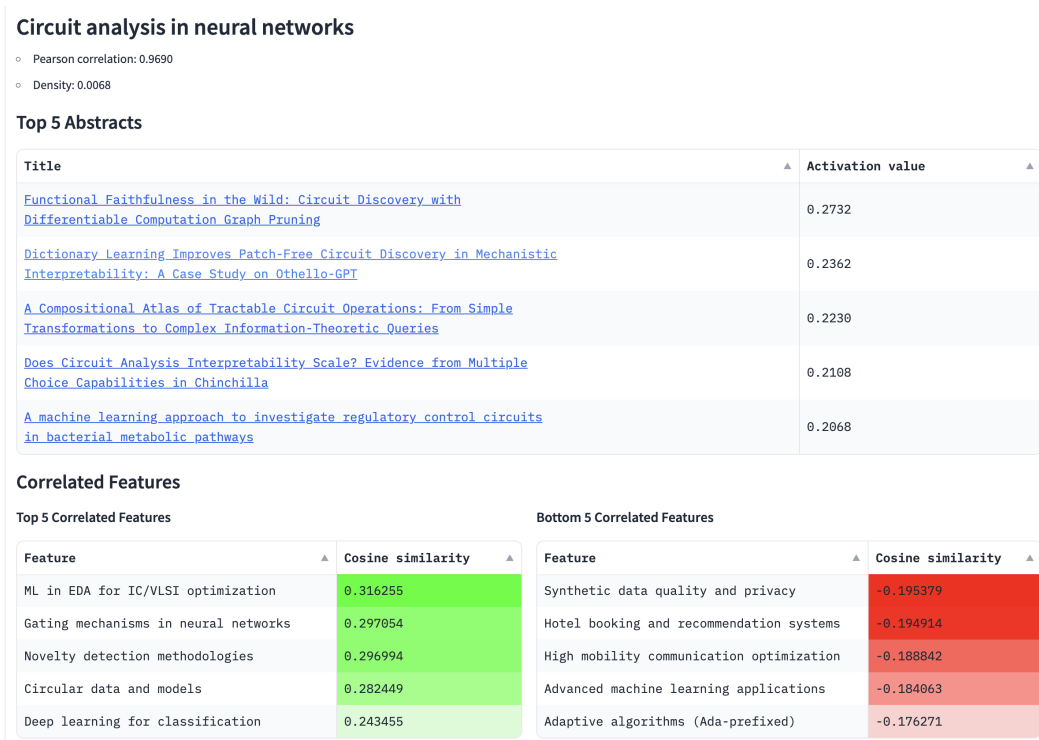


Figure 11: Individual feature visualisation for the “Circuit analysis in neural networks” feature in the computer science domain. The interface displays key interpretability metrics, top activating abstracts, correlated and co-occurring features, and an activation distribution histogram. Further information (not shown in the image) includes co-occurring features and activation distribution.

632 **B.2 Feature Visualisation Tab**

633 The Feature Visualisation tab is divided into two sub-tabs: Individual Features and Feature Families.  
 634 This section of the application directly relates to our analysis of SAE features and feature families  
 635 discussed in Sections 3 and 4.

636 **B.2.1 Individual Features**

637 For any selected feature, this tab displays:

- 638 • Top 5 activating abstracts, demonstrating the semantic content captured by the feature
- 639 • Top and bottom 5 correlated features, illustrating the relationships between different SAE  
 640 features
- 641 • Top 5 co-occurring features, showing which features tend to activate together
- 642 • A histogram of activation values, providing insight into the feature’s behavior across the  
 643 corpus
- 644 • The most similar features in SAE16 and SAE32

645 **B.2.2 Feature Families**

646 The Feature Families tab in our web application offers an in-depth exploration of related features  
 647 discovered by our sparse autoencoder. We show an example feature family in Figure 12.

648 The table displays the parent feature (superfeature) and its child features, along with key metrics,  
 649 such as the name of the parent and child features, the frequency of co-occurrence between the child  
 650 feature and the parent feature, ranging from 0 to 1, and the F1 Score and Pearson correlation.

651 The interactive directed graph provides a visual representation of the feature family structure. Each  
 652 node represents a feature. The size of the node corresponds to the feature’s density (frequency of

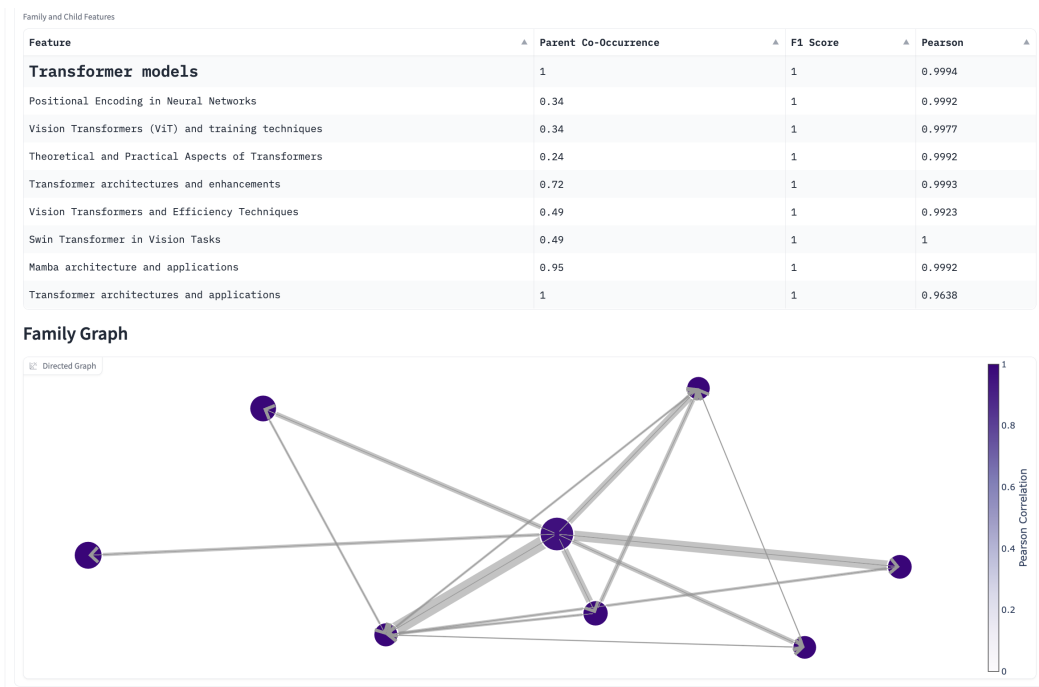


Figure 12: Directed graph visualization of a transformer models feature family. Nodes represent individual features, with size indicating feature density and color intensity showing Pearson correlation. Edges depict relationships between features, with arrow direction pointing from more general to more specific concepts. Users can hover over nodes to view detailed feature information.

653 activation), while the color intensity indicates the Pearson correlation (interpretability). Arrows  
 654 between nodes show relationships between features, with the direction typically pointing from more  
 655 general to more specific concepts. Users can hover over nodes to view detailed information about  
 656 each feature, including its name and log density.

## 657 C Automated interpretability details

### 658 C.1 Examples of features

659 We show some examples of perfectly interpretable features (Pearson correlation  $> 0.99$ ) in Table 3.  
 660 The strength of the activation of the feature on its top 3 activating abstracts is shown in parentheses  
 661 next to the abstract title.



<b>Feature</b>			
<b>Astronomy</b>			
Cosmic Microwave Background	CMB map-making and power spectrum estimation (0.1708)	How to calculate the CMB spectrum (0.1598)	CMB data analysis and sparsity (0.1581)
Periodicity in astronomical data	Generalized Lomb-Scargle analysis of decay rate measurements from the Physikalisch-Technische Bundesanstalt (0.1027)	Multicomponent power-density spectra of Kepler AGNs, an instrumental artefact or a physical origin? (0.0806)	RXTE observation of the X-ray burster 1E 1724-3045. I. Timing study of the persistent X-ray emission with the PCA (0.0758)
X-ray reflection spectra	X-ray reflection spectra from ionized slabs (0.3859)	The role of the reflection fraction in constraining black hole spin (0.3803)	Relativistic reflection: Review and recent developments in modeling (0.3698)
Critique or refutation of theories	What if string theory has no de Sitter vacua? (0.2917)	No evidence of mass segregation in massive young clusters (0.2051)	Ruling Out Initially Clustered Primordial Black Holes as Dark Matter (0.2029)
<b>Computer Science</b>			
Sparsity in Neural Networks	Two Sparsities Are Better Than One: Unlocking the Performance Benefits of Sparse-Sparse Networks (0.3807)	Truly Sparse Neural Networks at Scale (0.3714)	Topological Insights into Sparse Neural Networks (0.3689)
Gibbs Sampling and Variants	Herded Gibbs Sampling (0.2990)	Characterizing the Generalization Error of Gibbs Algorithm with Symmetrized KL information (0.2858)	A Framework for Neural Network Pruning Using Gibbs Distributions (0.2843)
Arithmetic operations in transformers	Arbitrary-Length Generalization for Addition in a Tiny Transformer (0.1828)	Carrying over algorithm in transformers (0.1803)	Understanding Addition in Transformers (0.1792)

Table 3: Activation strengths and titles for abstracts related to Astronomy and Computer Science features.

## 662 C.2 Automated interpretability prompts

663 We provide the prompts used for the Interpreter model and the Predictor model in the boxes below.  
664 Where this text is used, it represents an input to the model. We found that performance significantly  
665 increased when including the instruction to use ‘‘Occam’s razor’’, whereby the simplest feature at the  
666 appropriate level of granularity was selected.

### Interpreter Model Prompt

You are a meticulous <type> researcher conducting an important investigation into a certain neuron in a language model trained on <subject> papers. Your task is to figure out what sort of behaviour this neuron is responsible for – namely, on what general concepts, features, themes, methodologies or topics does this neuron fire? Here’s how you’ll complete the task:

**INPUT DESCRIPTION:** You will be given two inputs: 1) Max Activating Examples and 2) Zero Activating Examples.

1. You will be given several examples of text that activate the neuron, along with a number being how much it was activated. This means there is some feature, theme, methodology, topic or concept in this text that ‘excites’ this neuron.
2. You will also be given several examples of text that don’t activate the neuron. This means the feature, topic or concept is not present in these texts.

**OUTPUT DESCRIPTION:** Given the inputs provided, complete the following tasks.

1. Based on the MAX ACTIVATING EXAMPLES provided, write down potential topics, concepts, themes, methodologies and features that they share in common. These will need to be specific - remember, all of the text comes from subject, so these need to be highly specific subject concepts. You may need to look at different levels of granularity (i.e. subsets of a more general topic). List as many as you can think of. Give higher weight to concepts more present/prominent in examples with higher activations.
2. Based on the zero activating examples, rule out any of the topics/concepts/features listed above that are in the zero-activating examples. Systematically go through your list above.
3. Based on the above two steps, perform a thorough analysis of which feature, concept or topic, at what level of granularity, is likely to activate this neuron. Use Occam’s razor, as long as it fits the provided evidence. Be highly rational and analytical here.
4. Based on step 4, summarise this concept in 1-8 words, in the form FINAL: <explanation>. Do NOT return anything after these 1-8 words.

Here are the max-activating examples: <max activating examples>

Here are the zero-activating examples: <zero activating examples>

Work through the steps thoroughly and analytically to interpret our neuron.

667

### Predictor Model Prompt

You are a <subject> expert that is predicting which abstracts will activate a certain neuron in a language model trained on <subject> papers. Your task is to predict which of the following abstracts will activate the neuron the most. Here’s how you’ll complete the task:

**INPUT DESCRIPTION:** You will be given the description of the type of paper abstracts on which the neuron activates. This description will be short. You will then be given an abstract. Based on the concept of the abstract, you will predict whether the neuron will activate or not.

**OUTPUT DESCRIPTION:** Given the inputs provided, complete the following tasks.

1. Based on the description of the type of paper abstracts on which the neuron activates, reason step by step about whether the neuron will activate on this abstract or not. Be highly rational and analytical here. The abstract may not be clear cut - it may contain topics/concepts close to the neuron description, but not exact. In this case, reason thoroughly and use your best judgement. However, do not speculate on topics that are not present in the abstract.
2. Based on the above step, predict whether the neuron will activate on this abstract or not. If you predict it will activate, give a confidence score from 0 to 1 (i.e. 1 if you’re certain it will activate because it contains topics/concepts that match the description exactly, 0 if you’re highly uncertain). If you predict it will not activate, give a confidence score from -1 to 0.
3. Provide the final confidence score in the form PREDICTION: (your prediction) e.g. PREDICTION: 0.5. Do NOT return anything after this.

Here is the description/interpretation of the type of paper abstracts on which the neuron activates: <description>

Here is the abstract to predict: <abstract>

Work through the steps thoroughly and analytically to predict whether the neuron will activate on this abstract.

668

### 669 C.3 Exploring the effectiveness of smaller models

670 Although we eventually used gpt-4o-mini as the Predictor model, we initially did some ablations  
671 to understand how effective gpt-4o and gpt-3.5-turbo would be as different combinations of  
672 the Interpreter and Predictor models. We measured this by randomly sampling 50 features from  
673 our SAE64 (trained on astro-ph abstracts) and measuring the interpretability scores of different  
674 model combinations, in terms of both F1 score (does the model’s binary classification of a feature  
675 firing on an abstract agree with the ground-truth) and the Pearson correlation (described in the main  
676 body). Interestingly, we observe that using gpt-4o as the Interpreter and gpt-3.5-turbo as the  
677 Predictor leads to similar scores as using gpt-3.5-turbo for both, as shown in Figures 13 and  
678 Figures 14. This suggests that the challenging task in the autointerp is not necessarily labelling but  
679 rather predicting the activation of a feature on unseen abstracts.

680 Another observation is that using gpt-3.5-turbo as the Predictor only leads to a moderate degrada-  
681 tion of F1 score, it leads to a significant degradation of Pearson correlation. This is likely because  
682 we only use 6 abstracts for each feature prediction (3 positive, 3 negative) and thus there are only a  
683 few discrete F1 scores possible. Additionally, it appeared that gpt-3.5-turbo was generally less  
684 likely to assign higher confidence scores in either direction, with a much lower variance in assigned  
685 confidence than when gpt-4o was the Predictor. This affects Pearson correlation but not F1.

### 686 D Cross-domain features

687 The intersection between our cs.LG ( $n = 153, 146$ ) and astro.PH ( $n = 271, 492$ ) corpora contains  
688  $n = 330$  cross-posted papers. Motivated by these papers, as well as the observation of similar  
689 features re-occurring in models of different sizes (see Section 4), we search for the max cosine  
690 similarity feature between cs.LG and astro.PH SAEs at a fixed  $k$  and  $n_{dir}$ . As expected, we find

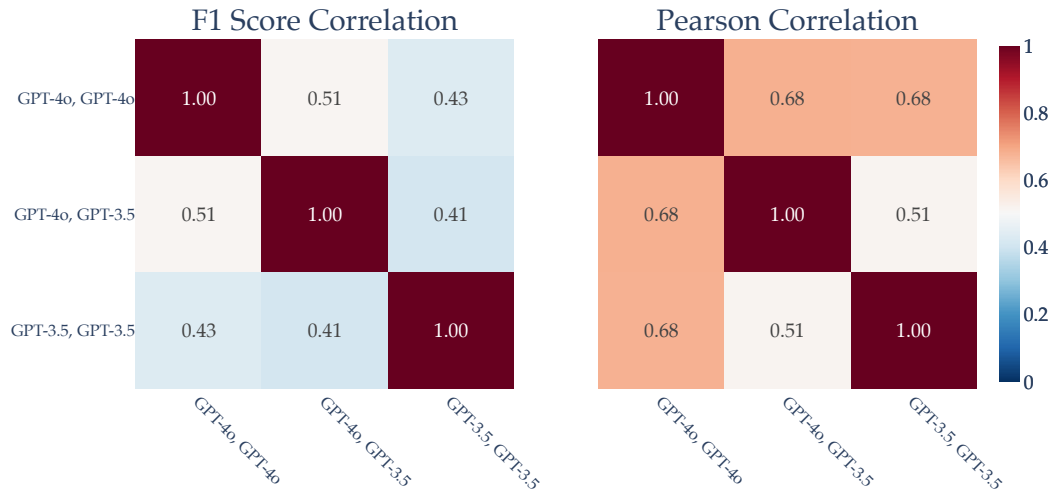


Figure 13: Correlation between F1 scores and Pearson correlation scores of different combinations of (labeller, predictor) models. Interestingly, using GPT-3.5 as the predictor appears to degrade performance similarly regardless of whether the feature was labelled by GPT-4o or GPT-3.5.

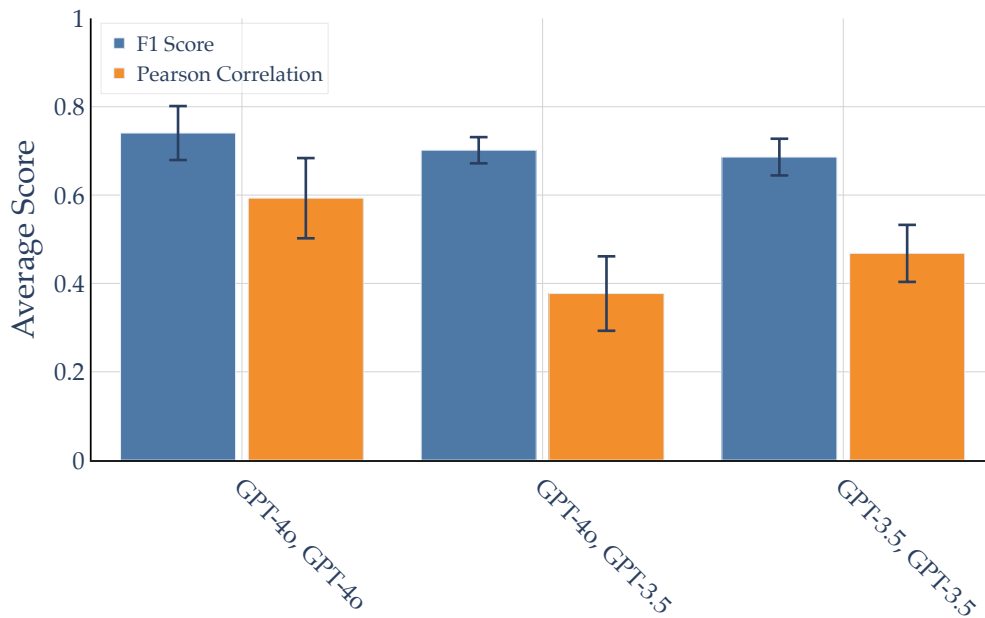


Figure 14: Mean F1 scores and Pearson correlations (according to ground-truth feature activations) across 50 randomly sampled features, for different combinations of (Interpreter, Predictor) models.

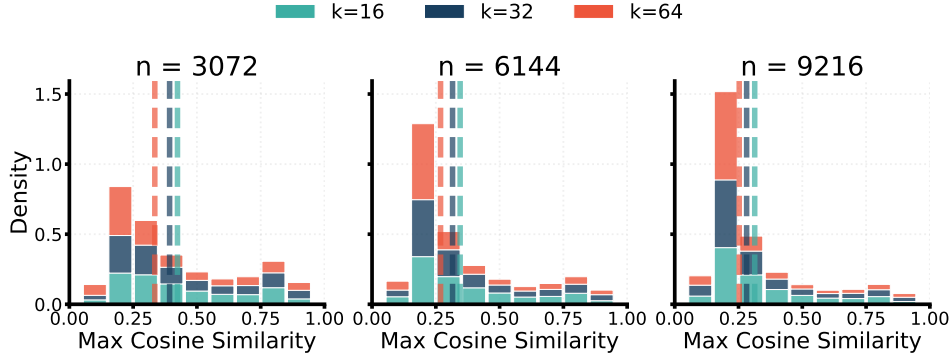


Figure 15: Maximum pair-wise cosine similarity of feature vectors between SAEs trained on different domains.

Feature Name (astro-ph)	Best Match (cs.LG)	Cosine Sim.	Activation Sim.	$\Delta F1$	$\Delta$ Pearson
Deep learning	CNNs and Applications	0.39	0.33	-0.2	-0.17
Generative Adversarial Networks	Generative Adversarial Networks (GANs)	0.61	0.26	0	0
Transformers	Transformer architectures and applications	0.5	0.33	0	-0
Artificial Neural Networks	Artificial Neural Networks (ANNs)	0.64	0.02	0	0
Artificial Intelligence	AI applications in diverse domains	0.61	0.45	0	0.02
Automation and Machine Learning	Automation in computational processes	0.9	0.77	-0.25	-0.47
Gaussian Processes	Gaussian Processes in Machine Learning	0.59	0.54	0	0.03
Regression analysis	Regression techniques and applications	0.81	0.53	0	-0.01

Table 4: Feature matches from the "Machine Learning" family (astroPH);  $k = 64$ ,  $n_{dir} = 9216$ .

691 significant mis-alignment between the vast majority of feature vectors between SAEs trained on  
 692 different domains, with mis-alignment increasing with  $k$  and  $n_{dir}$  (see Figure 15; this is unsurprising  
 693 given how  $k$  and  $n_{dirs}$  correlate with feature granularity).

694 However, a small subset of features appear in both sets of SAEs, with relatively high max cosine  
 695 similarity. For example, Table 4 shows the nearest cs.LG neighbours for every feature in the  
 696 astro.PH "Machine Learning" feature family (average cosine similarity = 0.59, average activation  
 697 similarity = 0.40). To test whether the features represent the same semantic concepts, we substitute the  
 698 natural language description of the best-match cs.LG feature for each listed astro.PH feature and  
 699 test the interpretability of the substituted descriptions; we find  $\Delta_{\text{pearson}} = -0.07$  and  $\Delta_{F1} = -0.06$ .  
 700 The existence of these features suggests that both sets of SAEs learn a semi-universal set of features  
 701 that span the domain overlap between astro.PH and cs.LG.

702 Interestingly, we find a number of near-perfectly aligned pairs (cosine similarity  $> 0.95$ ) of highly  
 703 interpretable features with little semantic overlap. A number of these features share similar wording  
 704 but not meaning, such as "Substructure in dark matter and galaxies" (astro-ph) and "Subgraphs and  
 705 their representations". Of these 10 feature pairs, the average activation similarity is 0.91.

## 706 E Feature family details

### 707 E.1 Feature splitting structures

708 Figure 16 shows an example of a recurrent feature across SAE sizes that does not exhibit feature  
 709 splitting. While the feature has extremely high activation and cosine similarity across every model  
 710 pair, each model only learns 1 feature in this direction. In Figures 17a and 17b we show two ex-  
 711 amples of feature splitting across SAE16 – SAE32 – SAE64 trained on astro-ph. 17a appears to  
 712 show canonical feature splitting as originally described in Bricken et al., 2023, with an increasing  
 713 number of features splitting the semantic space at each SAE size. There exists a top-level "period-  
 714 icity"/"periodicity detection" feature universal to all three SAEs, with relatively high similarity to  
 715 all other features, as well as novel, more granular features appearing in smaller SAEs, i.e. "Quasi-  
 716 periodic oscillations in blazars", which only appears in SAE64 and is highly dissimilar from other  
 717 split features.

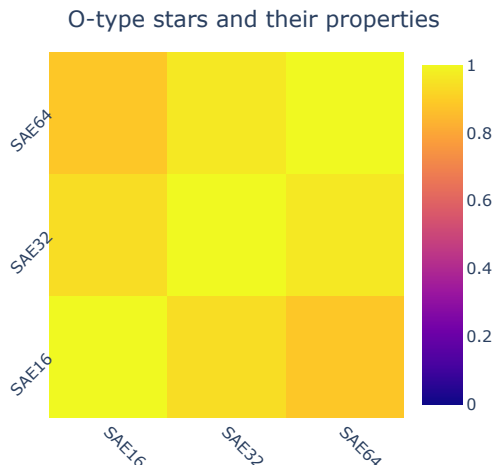


Figure 16: Recurrent features across SAEs trained on `astro-ph`; heatmap colored by activation similarity  $D$ ; all feature vector cosine similarities are  $> 0.98$ .

718 In contrast, 17b demonstrates nearest-neighbour features across models that do not exhibit semanti-  
 719 cally meaningful feature splitting. While the top-level “Luminous Blue Variables (LBVs)” feature  
 720 occurs at every model size, SAE64 also exhibits two additional features, “Lemaître-Tolman-Bondi  
 721 (LTB) Models” and “Lyman Break Galaxies (LBGs)”, that are highly dissimilar to each other, the  
 722 LBVs feature, and every other feature in the smaller models. We claim these are novel features,  
 723 occurring for the first time in SAE64, and that SAE16/SAE32 do not learn features for any related  
 724 higher-level concepts; instead, this grouping could be a spurious token-level correlation (LBV/LT-  
 725 B/LBG as similar acronyms).

726 **Feature triplets** In Figure 18a, we search for features that occur in  $n_{dirs} = 3072$  models and have  
 727 highly aligned features in larger ( $n_{dirs} = 6144, 9216$ ) models; we use this as a rough proxy for the  
 728 number of re-occurring features. We find that significantly more features re-occur between models  
 729 for higher  $k$ , with over 1100 feature triplets at  $> 0.95$  cosine similarity for  $k = 16$ ; as  $k$  increases,  
 730 the number of triplets drops sharply.

731 **Self-consistency** In 18b we show the set overlap between nearest-neighbour matches between  
 732 SAE16 and SAE64 found directly, and nearest-neighbour matches between SAE16 and SAE64 found  
 733 via nearest-neighbour matches to SAE32. If features exhibit perfectly clean splitting geometry, then  
 734 these two sets of SAE64 features should be consistent. However, we find that the distribution of set  
 735 overlap is roughly bimodal; other than triplet features with perfect overlap, overlap generally ranges  
 736 from 0 to 0.6. The vast majority of intersection = 1 sets are  $\leq 3$  features in size. This corroborates  
 737 findings in 9 which suggests features across models with different  $k$  are not well-aligned.

## 738 E.2 Feature family structure

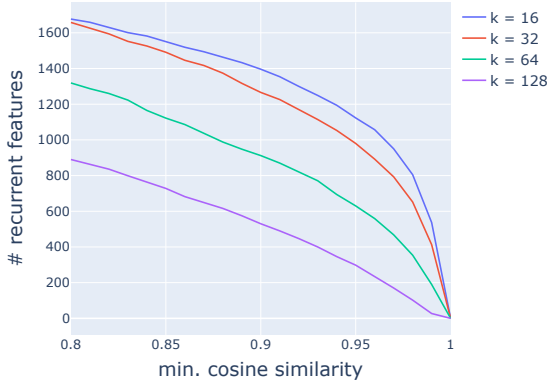
739 We compute feature family sizes (including the parent), co-occurrence ratios ( $\overline{R(p, C)}$ , see section 4),  
 740 and activation similarity ratios (computed identically to  $\overline{R(p, C)}$ , just using activation similarities).  
 741 Statistics for variants of `cs.LG` and `astro-ph` are shown in 19. We find a positive correlation  
 742 (Spearman = 0.22) between  $\overline{R(p, C)}$  and feature family interpretability.

743 We reproduce the projection method of Engels et al., 2024, running all documents through the SAE  
 744 and ablating features not in the feature family, to produce Figure 20. Visualizing the resulting principal  
 745 components confirms that the feature families we find do not represent manifolds or irreducible  
 746 multi-dimensional structures. We can instead think of feature families as linear subspaces in the  
 747 high-dimensional latent space; in fact, the component vectors can be seen in the lines of points  
 748 representing documents only activating on one feature in the family.

749 In 4 we use  $n = 3$  iterations of feature family construction. We select this hyper-parameter based off  
 750 Figure 21. In the first 2-3 iterations, removing parent nodes and re-constructing features preferentially

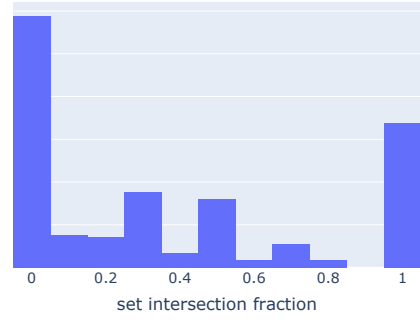






(a) Number of features from the smallest SAE that re-occur in all SAEs, by cosine similarity threshold.

Feature splitting (16-64 vs. 16-32-64)



(b) Overlap in the recovered SAE64 features, propagating nearest neighbors from SAE16-SAE64 vs. SAE16-SAE32-SAE64.

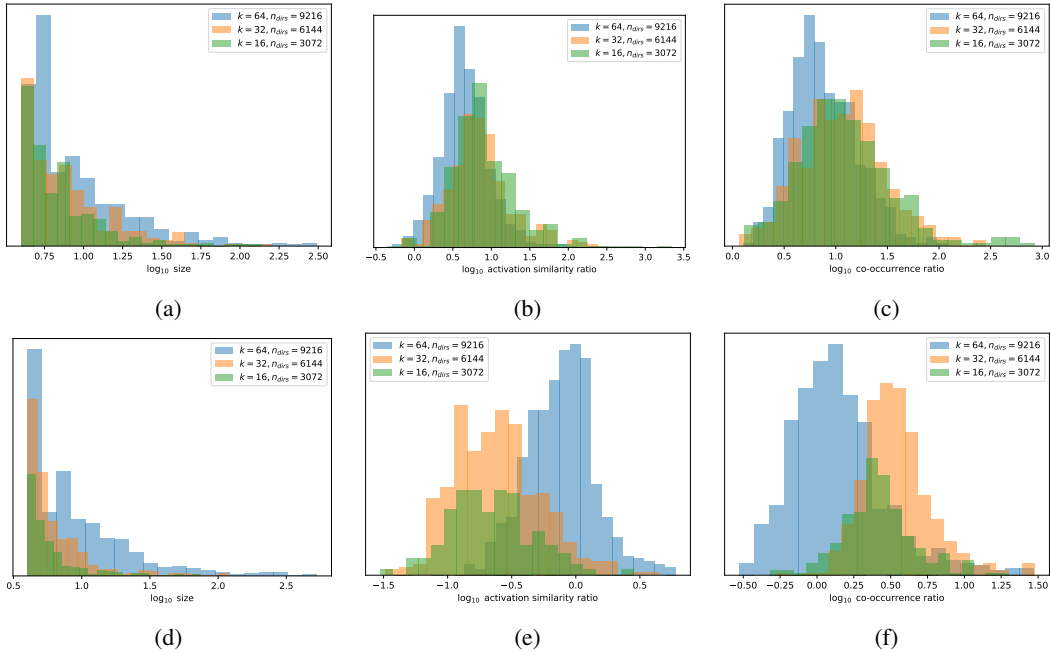


Figure 19: Feature families statistics (left: size; middle: activation similarity ratio; right: co-occurrence ratio,  $\overline{R(p, \mathcal{C})}$ );  $k = 64, n_{dir} = 9216$ .

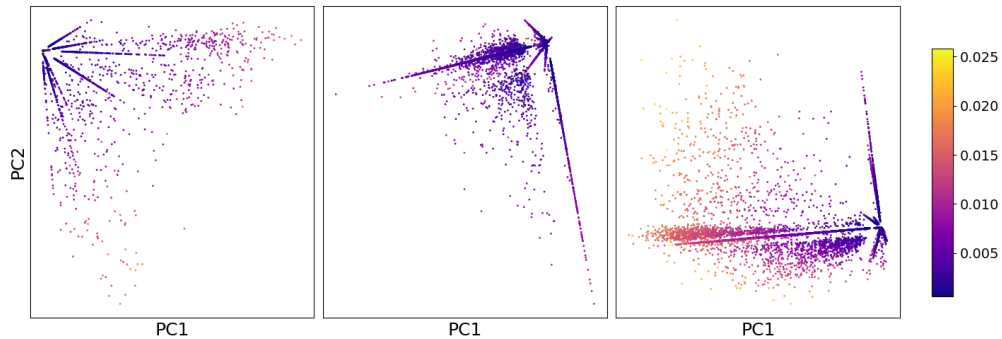


Figure 20: PCA projections of 3 example feature families from SAE64; points are latent representations of activating examples, colored by average activation for in-family features in the top  $k$ .

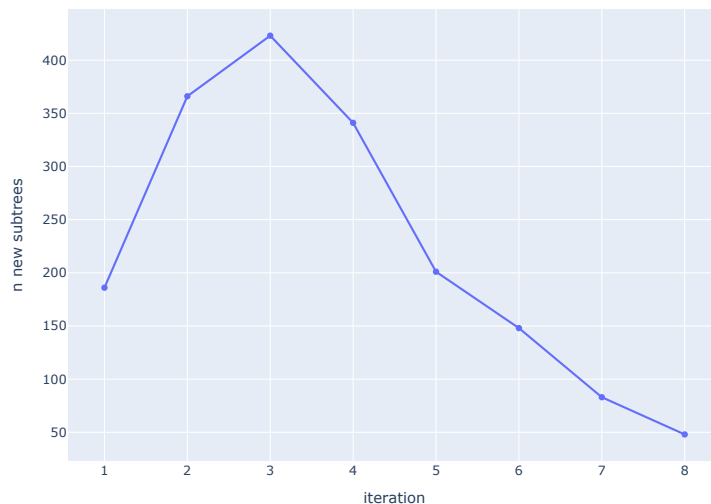


Figure 21: New feature families as a function of iteration; no deduplication is performed.

751 creates additional smaller families, suggesting iterations are necessary to fully explore the graph.  
 752 But given the sparse co-occurrences ( $C_{i,j} > 0.1$ ) used to build the graph, the number of additional  
 753 feature families found at each iteration drops off steeply after  $n = 3$ .

### 754 E.3 Feature family interpretability

755 We show example feature families and their interpretability scores in Figure 22.

## 756 F Exploring learned decoder weight matrices

757 **Encoder and decoder representations** Figure 23 reveals an intriguing relationship between feature  
 758 distinctiveness and the similarity of encoder and decoder representations in our sparse autoencoder.  
 759 In an ideal scenario with orthogonal features, encoder and decoder vectors would be identical, as the  
 760 optimal detection direction (encoder) would align perfectly with the representation direction (decoder).  
 761 This is because orthogonal features can be uniquely identified without interference. However, in our  
 762 high-dimensional space with more features than dimensions, perfect orthogonality is impossible due  
 763 to superposition.

764 The right panel of Figure 23 shows a negative correlation between a feature’s decoder-encoder cosine  
 765 similarity and its maximum similarity with other features. Features more orthogonal to others (lower  
 766 maximum similarity) tend to have more similar encoder and decoder representations. This aligns  
 767 with intuition: for more isolated features, the encoder’s detection direction can closely match the  
 768 decoder’s representation direction. Conversely, features with higher similarity to others require  
 769 the encoder to adopt a more differentiated detection strategy to minimise interference, resulting in  
 770 lower encoder-decoder similarity. The left panel, showing a mean cosine similarity of 0.57 between  
 771 corresponding encoder and decoder vectors, further emphasises this departure from orthogonality.  
 772 This phenomenon points to the importance of untied weights in sparse autoencoders.

773 **Clustering feature vectors** Motivated by structure in the feature activation graph, we explore whether  
 774 similar structure can be found in the decoder weight matrix  $W$  itself. Gao et al., 2024 find 2 such  
 775 clusters; we reproduce their method across our embeddings and SAEs, permuting the left singular  
 776 vectors  $U$  of  $W$  using a one-dimensional UMAP. We also experiment with permuting  $U$  and  $W$  using  
 777 reverse Cuthill-McKee. We do not find any meaningful block diagonal structure or clustering in  $W$ .

## 778 G Iterative encoding optimisation

779 We noted in Section 5 that intervening on a feature by up- or down-weighting its hidden representation  
 780 and then decoding is equivalent to directly adding the scaled feature vector to the final embedding.

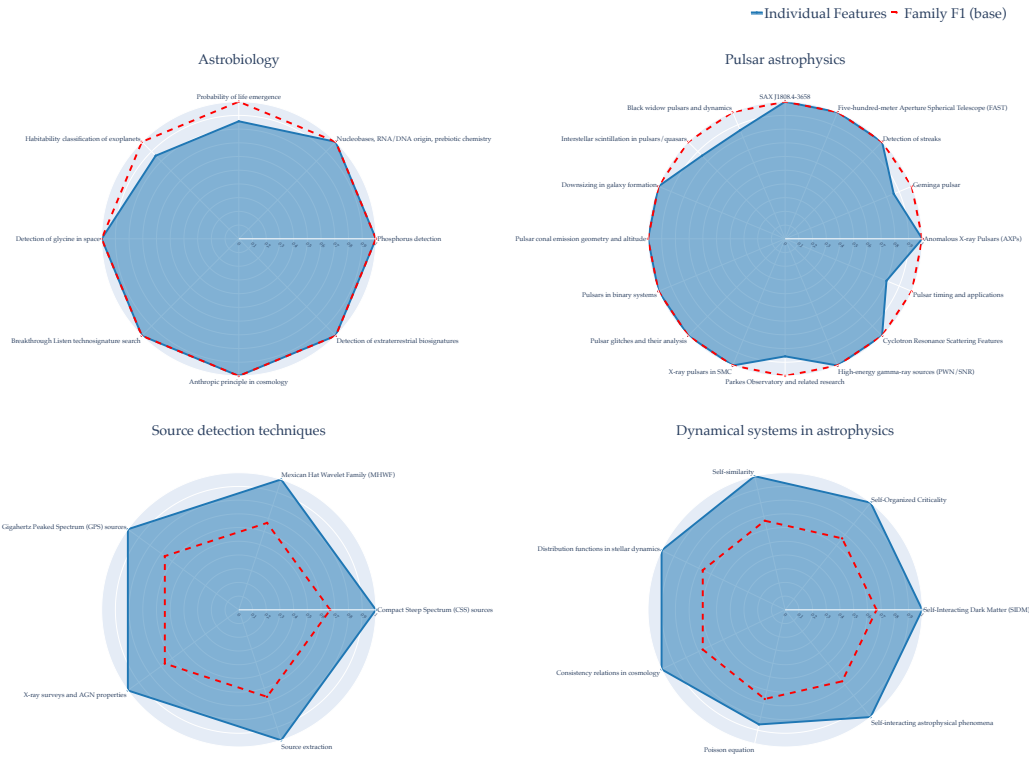


Figure 22: High-quality (top) and low-quality (bottom) feature families, scored through automated interpretability; radar charts show Pearson correlation scores for individual features (vertices) and the overall family (dashed line). While high-quality feature families truly have shared meaning, low-quality families appear to be mostly spurious and are not interpretable through short descriptions.

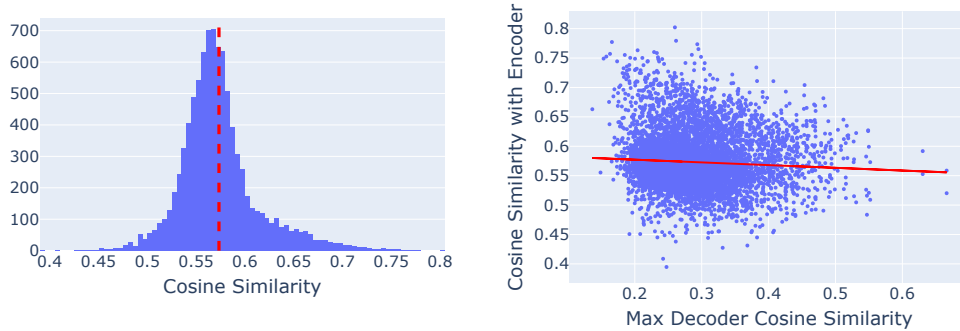


Figure 23: (Left) Cosine similarities between the encoder row and corresponding decoder column for SAE64 (cs.LG). The mean cosine similarity is 0.57, suggesting that encoder and decoder features are rather different, agreeing with Nanda (2023). (Right) We notice a slight negative correlation between a feature's decoder-encoder cosine similarity, and its maximum similarity with other features, possibly suggesting that features that are furthest removed from all other features in embedding space can have more similar corresponding decoders and encoder projections.



Figure 24: UMAP density plots along with LLM generated labels for SAE16 (left) and SAE64 (right) for the astro-ph features.

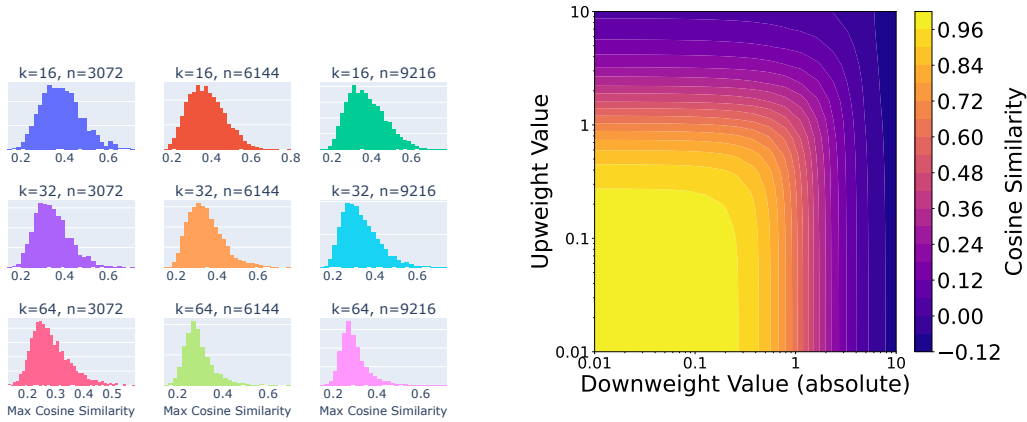


Figure 25: Distribution of maximum cosine similarity between a given feature vector and all other feature vectors, within the same SAE.

Figure 26: Cosine similarity between the original query embedding and the modified query embedding, with different values of upweighting random zero features and downweighting random active features.

781 To demonstrate this equivalence, let's consider an intervention on feature  $i$  by an amount  $\delta$ . The  
 782 modified hidden representation is  $\mathbf{h}' = \mathbf{h} + \delta \mathbf{e}_i$ , where  $\mathbf{e}_i$  is the  $i$ -th standard basis vector. Decoding  
 783 this modified representation gives  $\hat{\mathbf{x}}' = W_d \mathbf{h}' = W_d \mathbf{h} + \delta W_d \mathbf{e}_i = \hat{\mathbf{x}} + \delta \mathbf{w}_i$ , where  $\mathbf{w}_i$  is the  $i$ -th  
 784 column of  $W_d$ . Thus, intervening on the hidden representation and then decoding is equivalent to  
 785 directly adding the scaled feature vector to the original reconstruction.

786 We show in Figure 26 how cosine similarity between the original query embedding and the modified  
 787 query embedding changes as we change the upweighting and downweighting strength for different  
 788 features. Cosine similarity drops rapidly as soon as upweight or downweight exceeds 0.1.

789 There is an implicit challenge in SAE-based embedding interventions: the trade-off between steering  
 790 strength and precision. When directly manipulating feature activations, we observed that strong  
 791 interventions often led to unintended semantic shifts, activating correlated features and potentially  
 792 moving the embedding far from the SAE's learned manifold. Our goal is to achieve precise semantic  
 793 edits that express the desired feature strongly while minimising interference with unrelated features.  
 794 To this end, we developed an iterative optimisation approach that leverages the SAE's learned feature  
 795 space to find an optimal balance between these competing objectives.

796 Let  $\mathbf{x} \in \mathbb{R}^d$  be the original embedding,  $f_\theta(\cdot)$  the SAE encoder, and  $g_\phi(\cdot)$  the SAE decoder. We define  
 797 a target feature vector  $\mathbf{t} \in \mathbb{R}^k$  representing the desired feature activations after intervention, where  $k$   
 798 is the number of active features in our SAE. The iterative latent optimisation aims to find optimised  
 799 latents  $\mathbf{h}^*$  that satisfy:

$$\mathbf{h}^* = \operatorname{argmin}_{\mathbf{h}'} \{ \|f_\theta(g_\phi(\mathbf{h}')) - \mathbf{t}\|_2^2 \}$$

800 We solve this optimisation problem using gradient descent, starting from the initial latents  $\mathbf{h} = f_\theta(\mathbf{x})$   
 801 and iteratively updating  $\mathbf{h}'$ . We use the AdamW optimiser with a cosine annealing learning rate  
 802 schedule.

803 To evaluate the effectiveness of this approach, we compare it to a direct intervention method where we  
 804 simply set the target feature to a specific value in the latent space. For each abstract in our dataset, we  
 805 embed the abstract using an OpenAI embedding model to obtain  $\mathbf{x}$ . We then encode the embedding  
 806 to get initial latents  $\mathbf{h} = f_\theta(\mathbf{x})$ . We randomly select a target feature  $i$  and target value  $v$ . We then  
 807 apply both intervention methods: our iterative optimisation of  $\mathbf{h}'$  as described above, with  $\mathbf{t}_i = v$  and  
 808  $\mathbf{t}_j = \mathbf{h}_j$  for  $j \neq i$ , and direct intervention: setting  $\mathbf{h}'_i = v$  and  $\mathbf{h}'_j = \mathbf{h}_j$  for  $j \neq i$ .

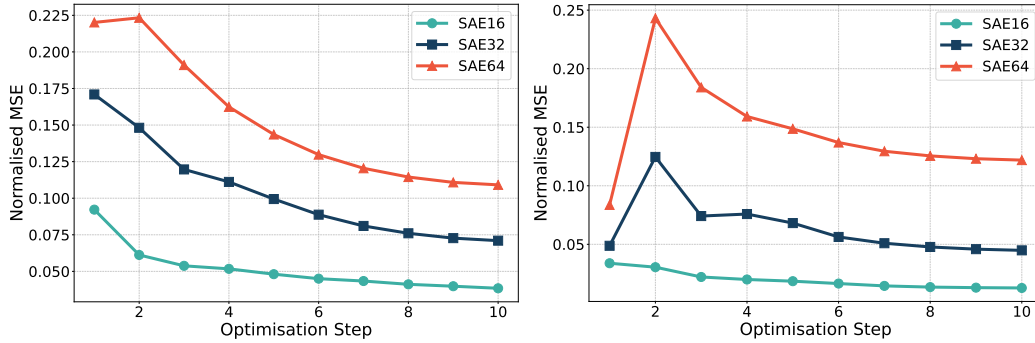


Figure 27: Normalised MSE at each of 10 steps across the iterative latent optimisation process. Left: Setting a random zero feature to active. Right: Setting a random active feature to zero.

809 Figure 27 (left panel) shows the trajectory of normalised MSE during the iterative optimisation process,  
 810 when setting a random zero feature to active. Similarly, the right panel shows the optimisation when  
 811 setting a random active feature to zero. Normalised MSE improves in the former case but not the  
 812 latter.