

---

# Interpretability as Compression: Reconsidering SAE Explanations of Neural Activations

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Sparse Autoencoders (SAEs) have emerged as a useful tool for interpreting the  
2 internal representations of neural networks. However, naively optimising SAEs for  
3 reconstruction loss and sparsity results in a preference for SAEs that are extremely  
4 wide and sparse. We present an information-theoretic framework for interpreting  
5 SAEs as lossy compression algorithms for communicating explanations of neural  
6 activations. We appeal to the Minimal Description Length (MDL) principle to  
7 motivate explanations of activations which are both accurate and concise. We  
8 further argue that interpretable SAEs require an additional property, “independent  
9 additivity”: features should be able to be understood separately. We demonstrate  
10 an example of applying our MDL-inspired framework by training SAEs on MNIST  
11 handwritten digits and find qualitatively more interpretable SAE features. We  
12 argue that using MDL rather than sparsity may avoid potential pitfalls with naively  
13 maximising sparsity such as undesirable feature splitting and that this framework  
14 naturally suggests new hierarchical SAE architectures which provide more concise  
15 explanations.

## 16 1 Introduction

17 Sparse Autoencoders (SAEs) (Le, 2013; Makhzani and Frey, 2013) were developed to learn a  
18 dictionary of sparsely activating features that describe a given dataset. They have recently become  
19 popular tools for interpreting the internal activations of large foundation language models, often  
20 finding human-understandable features (Sharkey et al., 2022; Huben et al., 2024; Bricken et al.,  
21 2023b). Researchers often use sparsity, the number of nonzero feature activations as measured by the  
22  $L_0$  norm, as a proxy for interpretability. SAEs are typically trained with an additional  $L_1$  penalty in  
23 their loss function to promote sparsity.

24 We adopt an information theoretic view of SAEs, inspired by Grünwald (2007), which views SAEs  
25 as explanatory tools that compress neural activations into communicable explanations. This view  
26 suggests that sparsity may appear as a special case of a larger objective: minimising the description  
27 length of the explanations. This operationalises Occam’s razor for selecting explanations: *all else*  
28 *equal, prefer the more concise explanation.*

## 29 2 SAEs are communicable explanations

30 SAEs aim to provide explanations of neural activations in terms of “features”<sup>1</sup>. Here we reformulate  
31 SAEs as solving a communication problem: suppose that we would like to transmit the neural  
32 activations  $x$  to a friend with some tolerance  $\varepsilon$ , either in terms of the reconstruction error or change in  
33 the downstream cross-entropy loss. Using the SAE as an encoding mechanism, we can approximate  
34 the representation of the activations in two parts. *First*, we send them the SAE encodings of the

35 activations  $z = Enc(x)$ . *Second*, we send them a decoder network  $Dec(\cdot)$  that recompiles these  
36 activations back to (some close approximation of) the neural activations,  $\hat{x} = Dec(z)$ .

37 This is closely analogous to *two-part coding schemes* (Grünwald, 2007) for transmitting a program  
38 via its source code and a compiler program that converts the source code into an executable format.  
39 Together the SAE activations and the decoder provide an **explanation** of the neural activations, based  
40 on the definition below.

41 **Definition 2.1** *An explanation  $e$  of some phenomena  $p$  is a statement  $e(p)$  for which knowing  $e(p)$   
42 gives some information about  $p$ . An explanation is typically a natural language statement<sup>2</sup>.*

43 The description length ( $DL$ ) of an explanation is the number of bits needed to transmit the explanation.  
44 For an SAE, this would be  $DL = |z|_{\text{bits}} + |Dec(\cdot)|_{\text{bits}}$ . The first term is  $O(n)$  and the second term is  
45  $O(1)$  in the dataset size so the first term dominates in the large data regime.

46 **Occam's Razor:** All else equal, an explanation  $e_1$  is preferred to explanation  $e_2$  if  $DL(e_1) < DL(e_2)$ .  
47 Intuitively, the simpler explanation is the better one. We can operationalise this as the Minimal  
48 Description Length (MDL) Principle for model selection: Choose the model with the shortest  
49 description length which solves the task. It has been observed that lower description length models  
50 often generalise better (MacKay, 2003).

51 **Definition 2.2** *We define the Minimal Description Length (MDL) as  $MDL_\epsilon(x) = \min DL(SAE)$   
52 where  $Loss(x, \hat{x}) < \epsilon$  and  $\hat{x} = SAE(x)$ . We say an SAE is  $\epsilon$ -MDL-optimal if it obtains this  
53 minimum.*

### 54 3 Interpretability requires independent additivity

55 Following Occam's razor we prefer simpler explanations, as measured by description length. But  
56 SAEs are not intended to simply give compressed explanations. They are also intended to give  
57 explanations that are interpretable and ideally human-understandable.

58 SAE features can be interpreted either as **causal results** of the model inputs (which we can see  
59 by analyzing feature activation patterns) or they can be interpreted as **causes** of the model outputs  
60 (which we can see through conducting interventions on the features and seeing the downstream  
61 effects). In both cases, we want to be able to understand each SAE feature independently, without  
62 needing to control for the activations of the other features. If all the feature activations are causally  
63 entangled—as is the case for the dense neural activations themselves—then they are not interpretable.  
64 Note that for  $D$  features there are  $O(D^2)$  pairs of features and  $\sum_i^K \binom{D}{i}$  possible sets of features  
65 which is much too large for humans to hold in working memory. So for feature explanations to be  
66 human-understandable we cannot have the all the features being entangled such that understanding a  
67 single concept requires understanding arbitrary feature interactions.

68 Hence, for interpretability, we need to be able to understand features independently of each other  
69 such that understanding a collection of features together is equivalent to understanding all the features  
70 separately. We call this property **independent additivity**, defined below.

71 **Definition 3.1** *Independent Additivity: An explanation  $e$  based on a vector of feature activations  
72  $\vec{z} = \sum_i \vec{z}_i$  is independently additive if  $e(\vec{z}) \approx \sum_i e(\vec{z}_i)$ . We say that a set of features  $z_i$  are  
73 independently additive if they can be understood independently of each other and the explanation of  
74 the sum of the features is the sum of the explanations of the features<sup>3</sup>.*

---

<sup>1</sup>Here we use the term "feature" as is common in the literature to refer to a linear direction which corresponds to a member of the set of a (typically overcomplete) basis for the activation space. Ideally the features are relatively monosemantic and correspond to a single (causally relevant) concept. We make no guarantees that the features found by an SAE are the "true" generating factors of the system.

<sup>2</sup>We will treat SAE activations and feature vectors as explanations themselves. Technically, we would want to do the additional step of interpreting their activation patterns or the results of causal interventions to get a natural language statement.

<sup>3</sup>Note that here the notion of summation depends on the explanation space. For natural language explanations, summation of adjectives is typically concatenation ("big" + "blue" + "bouncy" + "ball" = "The big blue bouncy ball"). For neural activations, summation is regular vector addition ( $\hat{x} = Dec(\vec{z}) = \sum_i Dec(z_i)$ ).

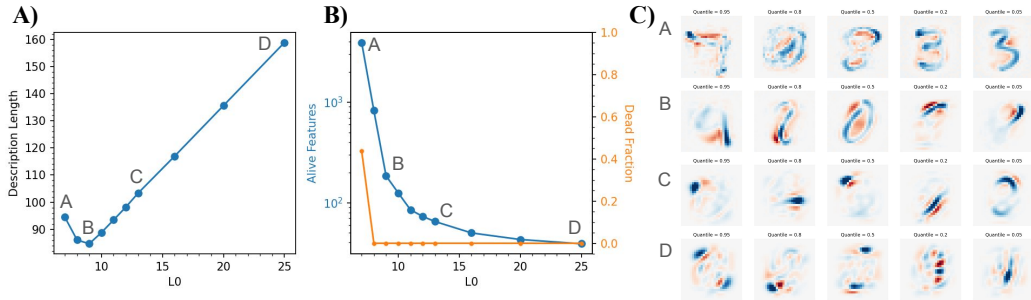


Figure 1: Finding the minimal description length (MDL) solution for SAEs trained on MNIST. A) Description length vs sparsity ( $L_0$ ) for a set of hyperparameters with the same reconstruction error. B) Plot of the number of alive features as a function of sparsity ( $L_0$ ). C) A random sample of SAE features at the 95th, 80th, 50th, 20th, and 5th percentiles of feature density respectively.

75 The independent additivity condition is directly analogous to the "composition as addition" property of  
 76 the Linear Representation Hypothesis (LRH) discussed in Olah (2024). *Independent additivity* relates  
 77 to the SAE features being composable via addition with respect to the explanation - this is a property  
 78 of the SAE Decoder. In the Linear Representation Hypothesis however, *Composition as Addition* is  
 79 about the underlying true features (i.e. the generating factors of the underlying distribution), which is  
 80 a property of the underlying distribution.

81 It is immediate from the definition that Independent Additivity holds for linear decoders however, we  
 82 note that this condition also allows for more general decoder architectures. For example, features can  
 83 be arranged to form a collection of directed trees, shown in fig. 3, where arrows represent the property  
 84 "the child node can only be active if the parent node is active"<sup>4</sup>. Here each feature still corresponds to  
 85 its own vector direction in the decoder. Since each child feature has a single path to its root feature,  
 86 there are no interactions to disentangle and the independent additivity property still holds, in that  
 87 each *tree* can be understood independently in a way that's natural for humans to understand, as a  
 88 multi-dimensional feature. An advantage of the directed-tree SAE decoder structure is that it can be  
 89 more description-length efficient as shown in fig. 5.

## 90 4 The MDL-SAE finds interpretable and composable features for MNIST

91 To achieve the same loss, higher sparsity (lower  $L_0$ ) typically requires a larger dictionary, so there's  
 92 an inherent trade-off between decreasing  $L_0$  and decreasing the dictionary size in order to reduce  
 93 description length. We explore this trade-off in MNIST below and for the GPT-2 language model in  
 94 appendix B.

95 Lee (2001) describe the classical method for using the Minimal Description Length (MDL) criteria  
 96 for model selection. Here we choose between model hyperparameters (in particular the SAE width  
 97 and expected  $L_0$ ) for the optimal SAE. Our algorithm for finding the MDL-SAE solution and details  
 98 for this case study are given in appendix F. We trained SAEs on the MNIST dataset of handwritten  
 99 digits (LeCun et al., 1998) and find the set of hyperparameters resulting in the same test MSE. We  
 100 see three basic regimes:

- 101 • **High  $L_0$ , narrow SAE width** (C, D in fig. 1): Here, the description length (DL) is linear  
 102 with  $L_0$ , suggesting that the DL is dominated by the number of bits needed to represent the  
 103  $L_0$  nonzero floats. The features appear as small sections of digits that could be relevant to  
 104 many digits (C) or start to look like dense features that one might obtain from PCA (D).
- 105 • **Low  $L_0$ , wide SAE width** (A in fig. 1): Though  $L_0$  is small, the DL is large because as the  
 106 SAE becomes wider, additional bits are required to specify which activations are nonzero.  
 107 The features appear closer to being full digits, i.e. similar to samples from the dataset.

<sup>4</sup>In practice, we typically expect feature trees to be shallow structures which capture causal relationships between highly related features. A particularly interesting example of this structure is a group-sparse autoencoder where linear subspaces are densely activated together.

108 • **The MDL solution** (B in fig. 1): There’s a balance between the two contributions to the  
109 description length. The features appear like longer line segments or strokes for digits, but  
110 could apply to multiple digits.

111 In this example, the MDL solution finds a meaningful decomposition of digits into stroke-like features.  
112 More dense SAEs find less interpretable point-like features, while sparser SAEs find features that  
113 resemble examples from the dataset and fail to decompose the digits into reusable and composable  
114 features.

## 115 5 Optimising for MDL can reduce undesirable feature splitting

116 In large language models, SAEs with larger dictionaries learn finer-grained versions of features  
117 learned in smaller SAEs, a phenomenon known as "feature splitting" (Bricken et al., 2023b). Feature  
118 splitting that introduces a novel conceptual distinction is desirable but some feature splitting—for  
119 example, learning dozens of features representing the letter "P" in different contexts (Bricken et al.,  
120 2023b)—is undesirable and can waste dictionary capacity while not giving more explanatory power.

121 A toy model of undesirable feature splitting is an SAE that represents the AND of two boolean  
122 features,  $A$  and  $B$ , as a third feature direction. The two booleans represent whether the feature vectors  
123  $v_A$  and  $v_B$  are present or not, so there are four possible activations: 0,  $v_A$ ,  $v_B$ , and  $v_A + v_B$ . See  
124 appendix E for details on our model of feature splitting.

125 Even though feature splitting always results in a lower  $L_0$ , it does not always result in the smallest  
126 description length. The phase diagram in fig. 4 shows the case where  $p_A = p_B$ . If the correlation  
127 coefficient  $\rho$  between  $A$  and  $B$  is small then representing only  $A$  and  $B$ , but not  $A \wedge B$ , takes fewer  
128 bits so the preferred solution avoids feature splitting. However, if the correlation is large, then feature  
129 splitting is preferred since  $A \wedge B$  occurs frequently enough that explicitly representing it reduces the  
130 description length. In this way, minimizing description length can limit the amount of undesirable  
131 feature splitting and gives us a concrete decision criteria to understand when we might expect feature  
132 splitting.

## 133 6 Related Work

134 Bricken et al. (2023a) also consider how information measures relate to SAEs and find that "bounces"  
135 in entropy correspond to dictionary sizes with the correct number of features in synthetic experiments.  
136 We find a similar bounce in description length in a non-synthetic experiment. We go further by  
137 studying several examples where minimal description length gives more intuitive features and discuss  
138 more description-efficient SAE architectures.

139 As in Ramirez and Sapiro (2012), we use the MDL approach for the Model Selection Problem  
140 using the criteria that the best model for the data is the model that captures the most useful structure  
141 from the data. Chan et al. (2024) use Mechanistic Interpretability techniques to generate compact  
142 formal guarantees (i.e. proofs) of model performance and also note a deep connection between  
143 interpretability and compression.

## 144 7 Conclusion

145 In this work, we have presented an information-theoretic perspective on Sparse Autoencoders as  
146 explainers for neural network activations. Using the MDL principle, we provide some theoretical  
147 motivation for existing SAE architectures and hyperparameters. We also hypothesise a mechanism  
148 for, and criteria to describe, the commonly observed phenomena of feature splitting. In the cases  
149 where feature splitting can be seen as undesirable for downstream applications, we hope that, using  
150 this theoretical framework, the prevalence of undesirable feature splitting could be decreased in  
151 practical modelling settings.

## 152 8 Debunking Challenge

### 153 8.1 What commonly-held position or belief are you challenging?

154 In Mechanistic Interpretability, a commonly held belief is that interpretable explanations consist of  
155 sparse latents and in particular that sparsity is an operationalisable proxy for interpretability that we  
156 can use both in our loss functions and for model selection (Sharkey et al., 2022; Huben et al., 2024;  
157 Bricken et al., 2023b; Olah, 2024; Gao et al., 2024).

158 In this way, researchers typically try to attain SAEs which perform well on a (reconstruction error,  
159 sparsity)-Pareto frontier. That is, they seek to reconstruct the data effectively using as few SAE  
160 latent features as possible, trading off reconstruction error and sparsity. Though this framing has  
161 been produced useful features in an unsupervised manner in some cases, we note that optimising for  
162 sparsity has many undesirable outcomes:

- 163 • **Feature Splitting:** In section 5 and appendix E, we give a model of undesirable feature  
164 splitting in which an SAE learns to represent the AND of two genuine model features as a  
165 third latent feature direction. If the model has sufficient width, then, from the perspective of  
166 sparsity, it is always beneficial to represent ANDs of features which co-occur even once  
167 in the dataset. We can see that this leads to unrestrained, undesirable feature splitting as  
168 several composite features get their own direction in the SAE even if they are not salient  
169 features either to the human interpreters or to the model.
- 170 • **Strange Limiting Properties for sparsity:** Suppose that we're jointly optimising for  
171 sparsity and reconstruction error with the SAE width as a free hyperparameter. How should  
172 we expect the width to change as we optimise? It is immediate to see that the width should  
173 grow to be extremely large: the natural solution to this optimisation problem is to take  
174 the width of the SAE to be equal to the number of possible neural activation inputs (i.e.  
175  $D = (\text{vocab size})^{\text{seq len}}$ ). In this case, we have the sparsity given as  $L_0 = 1$  and MSE  
176 reconstruction error as 0. Though this solution is optimal given the problem statement, it is  
177 difficult to see this solution as a valuable tool for interpreting the neural activations.

### 178 8.2 How are your results in tension with this commonly-held position?

179 We argue from the information-theoretic perspective of viewing SAEs as explanatory tools that  
180 compress neural activations into communicable explanation. From this perspective, minimising  
181 sparsity appears not as the true optimisation goal but rather as a proxy for minimising description  
182 length (i.e. conciseness).

183 Under the MDL paradigm, we instead are able to overcome the two previously presented issues. For  
184 feature splitting, fig. 4 shows that MDL SAEs have a clear decision boundary which describes which  
185 feature splitting is deemed effective and so naturally reduce the prevalence of undesirable feature  
186 splitting in our feature dictionary. This results in SAEs which are subjectively more interpretable and  
187 speculatively appear to be more aligned with the model's computation.

### 188 8.3 How do you expect your submission to affect future work?

189 We expect future work to optimise for the (reconstruct error, description length)-Pareto frontier rather  
190 than the (reconstruction error, sparsity)-Pareto frontier. We show that this approach naturally suggests  
191 SAE architectures which admit more human-interpretable features. Promising architectures for future  
192 work include hierarchical and group-sparse SAEs.

193 Our work also naturally suggests a different approach to the unbounded search for ever-wider SAEs  
194 as present in Templeton et al. (2024). We expect future work to focus more on obtaining disentangled,  
195 causally relevant, interpretable features rather than pushing on the size of the dictionary.

196 In particular, though (Engels et al., 2024) suggest that not all language model model representations  
197 are 1-d subspaces, we note that it is hard to successfully use this fact to build better SAE architectures.  
198 This is because, from a sparsity perspective, it is still better to instead use feature splitting to find  
199 linear directions rather than actually taking advantage of the inherent geometry of the feature space.  
200 With MDL-SAEs, it becomes feasible to successfully use the geometry of feature space to reduce the  
201 description length of the explanation (at the expense of sparsity), giving more interpretable features.

## 202 References

- 203 J. Bloom. Open source sparse autoencoders for all residual stream layers of gpt2-small. *AI Alignment*  
204 *Forum*, 2024. URL [https://www.alignmentforum.org/posts/f9EgflSurAiqRjySD/](https://www.alignmentforum.org/posts/f9EgflSurAiqRjySD/open-source-sparse-autoencoders-for-all-residual-stream)  
205 [open-source-sparse-autoencoders-for-all-residual-stream](https://www.alignmentforum.org/posts/f9EgflSurAiqRjySD/open-source-sparse-autoencoders-for-all-residual-stream).
- 206 T. Bricken, J. Batson, A. Templeton, A. Jermyn, T. Henighan, and C. Olah. Features as the simplest  
207 factorization. *Transformer Circuits Thread*, 2023a. URL [https://transformer-circuits.](https://transformer-circuits.pub/2023/may-update/index.html#simple-factorization)  
208 [pub/2023/may-update/index.html#simple-factorization](https://transformer-circuits.pub/2023/may-update/index.html#simple-factorization).
- 209 T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner, C. Anil, C. Denison,  
210 A. Askell, R. Lasenby, Y. Wu, S. Kravec, N. Schiefer, T. Maxwell, N. Joseph, Z. Hatfield-  
211 Dodds, A. Tamkin, K. Nguyen, B. McLean, J. E. Burke, T. Hume, S. Carter, T. Henighan,  
212 and C. Olah. Towards monosemanticity: Decomposing language models with dictionary learn-  
213 ing. *Transformer Circuits Thread*, 2023b. URL [https://transformer-circuits.pub/2023/](https://transformer-circuits.pub/2023/monosemantic-features/index.html)  
214 [monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 215 B. Bussmann, P. Leask, and N. Nanda. Batchtopk: A simple improvement for topk-saes. *AI Alignment*  
216 *Forum*, 2024. URL [https://www.alignmentforum.org/posts/Nkx6yWZNbAsfvic98/](https://www.alignmentforum.org/posts/Nkx6yWZNbAsfvic98/batchtopk-a-simple-improvement-for-topk-saes)  
217 [batchtopk-a-simple-improvement-for-topk-saes](https://www.alignmentforum.org/posts/Nkx6yWZNbAsfvic98/batchtopk-a-simple-improvement-for-topk-saes).
- 218 L. Chan, R. Agrawal, A. Garriga-Alonso, and J. Gross. Compact proofs of model performance via  
219 mechanistic interpretability. *AI Alignment Forum*, 2024. URL [https://www.alignmentforum.](https://www.alignmentforum.org/posts/bRsKimQcPTX3tNNJZ)  
220 [org/posts/bRsKimQcPTX3tNNJZ](https://www.alignmentforum.org/posts/bRsKimQcPTX3tNNJZ).
- 221 J. Engels, I. Liao, E. J. Michaud, W. Gurnee, and M. Tegmark. Not all language model features are  
222 linear, 2024. URL <https://arxiv.org/abs/2405.14860>.
- 223 L. Gao, T. D. la Tour, H. Tillman, G. Goh, R. Troll, A. Radford, I. Sutskever, J. Leike, and J. Wu.  
224 Scaling and evaluating sparse autoencoders. *arXiv preprint arXiv:2406.04093*, 2024.
- 225 P. D. Grünwald. *The minimum description length principle*. MIT press, 2007.
- 226 R. Huben, H. Cunningham, L. R. Smith, A. Ewart, and L. Sharkey. Sparse autoencoders find highly  
227 interpretable features in language models. In *The Twelfth International Conference on Learning*  
228 *Representations*, 2024. URL <https://openreview.net/forum?id=F76bwRSLek>.
- 229 Q. V. Le. Building high-level features using large scale unsupervised learning. In *2013 IEEE*  
230 *international conference on acoustics, speech and signal processing*, pages 8595–8598. IEEE,  
231 2013.
- 232 Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document  
233 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 234 T. C. Lee. An introduction to coding theory and the two-part minimum description length principle.  
235 *International statistical review*, 69(2):169–183, 2001.
- 236 D. J. MacKay. *Information theory, inference and learning algorithms*. Cambridge university press,  
237 2003.
- 238 A. Makhzani and B. Frey. K-sparse autoencoders. *arXiv preprint arXiv:1312.5663*, 2013.
- 239 C. Olah. Circuits updates - july 2024, linear representations. *Transformer Circuits Thread*, 2024.  
240 <https://transformer-circuits.pub/2024/july-update/index.html#linear-representations>.
- 241 I. Ramirez and G. Sapiro. An mdl framework for sparse coding and dictionary learning. *IEEE*  
242 *Transactions on Signal Processing*, 60(6):2913–2927, 2012.
- 243 D. Salomon. *Variable-length codes for data compression*. Springer Science & Business Media, 2007.
- 244 L. Sharkey, D. Braun, and B. Millidge. Taking features out of su-  
245 perposition with sparse autoencoders. *AI Alignment Forum*, 2022.  
246 URL [https://www.alignmentforum.org/posts/z6QQJbtPkEAX3Aojj/](https://www.alignmentforum.org/posts/z6QQJbtPkEAX3Aojj/interim-research-report-taking-features-out-of-superposition)  
247 [interim-research-report-taking-features-out-of-superposition](https://www.alignmentforum.org/posts/z6QQJbtPkEAX3Aojj/interim-research-report-taking-features-out-of-superposition).

248 A. Templeton, T. Conerly, J. Marcus, J. Lindsey, T. Bricken, B. Chen, A. Pearce, C. Citro, E. Ameisen,  
249 A. Jones, H. Cunningham, N. L. Turner, C. McDougall, M. MacDiarmid, C. D. Freeman, T. R.  
250 Summers, E. Rees, J. Batson, A. Jermyn, S. Carter, C. Olah, and T. Henighan. Scaling monoseman-  
251 ticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*,  
252 2024. URL [https://transformer-circuits.pub/2024/scaling-monosemanticity/](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html)  
253 [index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).

254 **A SAE communication protocol**

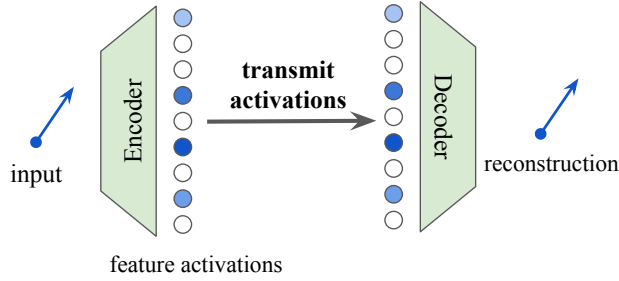


Figure 2: A schematic showing a sparse autoencoder (SAE) being used to communicate an input by transmitting the encoded activations and decoding them into a reconstruction of the input.

255 **B SAEs should be sparse, but not too sparse**

256 Naively we might see SAEs as decompressing neural activations which contain densely packed  
 257 features in superposition. To see that SAEs are producing compressed explanations of activations we  
 258 must note that the inherent feature sparsity means that it is more efficient to communicate SAE latent  
 259 features rather than neural activations even though the dimension of the latent dimension is higher.

260 The description length for a set of SAE activations (under independent additivity) with distribution  
 261  $p(z)$  is given by  $H(p) = \sum_{z \in Z} -p(z) \log_2 p(z)$ . For exposition, consider a simpler formulation  
 262 where we directly consider the bits needed without prior knowledge of the distributions. For a set of  
 263 feature activations with  $L_0$  nonzero elements out of  $D$  dictionary features, an upper bound on the  
 264 description length is

$$DL \lesssim L_0(B + \log_2 D) \tag{1}$$

265 where  $B$  is the effective precision of each float and  $\log_2 D$  is the number of bits required to specify  
 266 which features are active. To achieve the same loss, higher sparsity (lower  $L_0$ ) typically requires a  
 267 larger dictionary, so there’s an inherent trade-off between decreasing  $L_0$  and decreasing the dictionary  
 268 size in order to reduce description length.

269 As an illustrative example, in Appendix B, we compare reasonable hyperparameters for GPT-2 SAEs  
 270 to dense/narrow and sparse/wide extreme hyperparameters. We show that an SAE (Bloom, 2024)  
 271 has a description length of approximately 1405 bits per input token, compared to 5376 bits for  
 272 transmitting the dense neural activations and 13,993 bits for a one-hot encoding of all possible token  
 273 sequences of length 128. Here the SAE at intermediate sparsity and width has the lower description  
 274 length.

275 **C Independently Additive SAE Architectures**

276 We show examples of different SAE architectures that satisfy independent additivity in fig. 3

277 **D Comparison of GPT-2 SAE hyperparameters**

- 278 • **Reasonable SAEs:** Bloom (2024)’s open-source SAEs for GPT-2 layer 8 have  $L_0 = 65$ ,  
 279  $D = 25,000$ . Given  $B = 7$  bits per nonzero float (8-bit quantization with the sign fixed to  
 280 positive), the description length per input token is 1405 bits.
- 281 • **Dense Activations:** A dense representation that still satisfies independent additivity would  
 282 be to send the neural activations directly instead of training an SAE. GPT-2 has a model size  
 283 of  $d = 768$ , the description length is simply  $DL = B d = 5376$  bits per token.
- 284 • **One-hot encodings:** At the sparse extreme, our dictionary has a row for each neural  
 285 activation in the dataset, so  $L_0 = 1$  and  $D = (\text{vocab size})^{\text{seq len}}$ . GPT-2 has a vocab size of



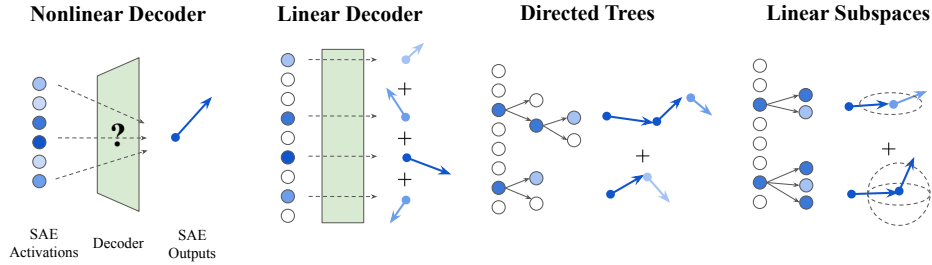


Figure 3: Examples of different SAE architectures. All but nonlinear decoders are compatible with independent additivity as feature activations correspond to adding a separate vector to the output. Architectures with directed tree decoders or which allow for vectors lying within a subspace are potentially more communication efficient since a child node can only be active if its parent node is active.

286 50,257 and the SAEs are trained 128 token sequences. All together this gives  $DL = 13,993$   
 287 bits per token.

288 Although the comparison is slightly unfair because the SAE is lossy (93% variance explained) and the  
 289 other cases are lossless, these calculations demonstrate that reasonable SAEs are indeed compressed  
 290 compared to the dense and sparse extremes. We hypothesise that the reason that we’re able to get this  
 291 helpful compression is that the true features from the generating process are themselves sparse.

292 Note the difference here from choosing models based on the reconstruction loss vs sparsity ( $L_0$ )  
 293 Pareto frontier. When minimising  $L_0$ , we are encouraging decreasing  $L_0$  and increasing  $D$  until  
 294  $L_0 = 1$ . Under the MDL model selection paradigm we are typically able to discount trivial solutions  
 295 like a one-hot encoding of the input activations and other extremely sparse solutions which make the  
 296 reconstruction algorithm analogous to a k-Nearest Neighbour classifier.

## 297 E Toy Model of Feature Splitting

298 **No Feature Splitting:** Say that the SAE only learns two boolean feature vectors,  $v_A$  and  $v_B$ , as  
 299 shown in fig. 4. It is still capable of reconstructing  $A \wedge B$  as the sum  $v_A + v_B$ . The  $L_0$  would simply  
 300 be the expectation of the boolean activations, so  $L_0 = p_A + p_B$  and the description length would be  
 301  $DL = H(p_A) + H(p_B)$  where  $H(p)$  is the entropy of a Bernoulli variable with probability  $p$ .

302 **Feature Splitting:** In this case, the SAE learns three mutually exclusive features.  $A \wedge B$  is explicitly  
 303 represented with the vector  $v_A + v_B$  while the two other features represent  $A \wedge \neg B$  and  $B \wedge \neg A$   
 304 with vectors  $v_A$  and  $v_B$ . This setup has the same reconstruction error but has lower  $L_0 = p_{A \wedge \neg B} +$   
 305  $p_{B \wedge \neg A} + p_{A \wedge B} = p_A + p_B - p_{A \wedge B}$  since the probabilities for  $A \wedge \neg B$ , say, are reduced as  
 306  $p_{A \wedge \neg B} = p_A - p_{A \wedge B}$ . Note that the  $L_0$  (sparsity) is necessarily lower than in the non-feature  
 307 splitting case.

## 308 F Details on determining the MDL-SAE

### 309 F.1 Algorithm

- 310 1. **Specify a tolerance level,  $\varepsilon$ , for the loss function.** The tolerance  $\varepsilon$  is the maximum allowed  
 311 value for the loss, either the reconstruction loss (MSE for the SAE) or the model’s cross-  
 312 entropy loss when intervening on the model to swap in the SAE reconstructions in place of  
 313 the clean activations. For small datasets using a reconstruction, the test loss should be used.
- 314 2. **Train a set of SAEs within the loss tolerance.** It may be possible to simplify this task by  
 315 allowing the sparsity parameter to also be learned.
- 316 3. **Find the effective precision needed for floats.** The description length depends on the float  
 317 quantisation. We typically reduce the float precision until the change in loss results in the  
 318 reconstruction tolerance level is exceeded.

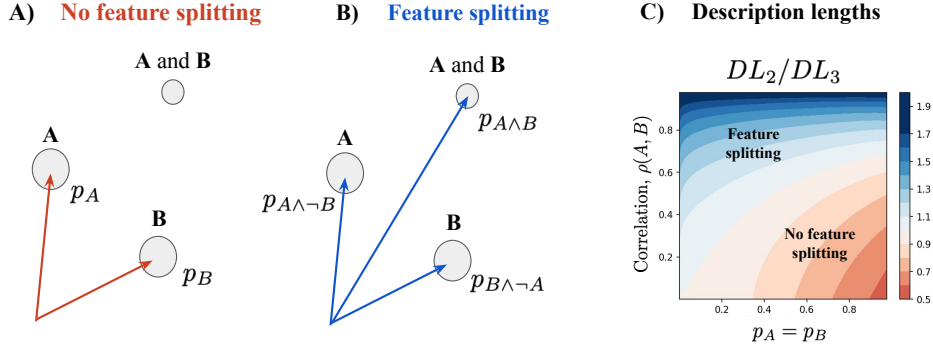


Figure 4: A toy model of undesirable feature splitting. The SAE can learn two boolean features without feature splitting (A) or three mutually exclusive boolean features with feature splitting (B) which always has lower  $L_0$ . Minimizing description length provides a decision boundary (C) for when feature splitting is preferred or not.

4. **Calculate description lengths.** With the quantised latent activations, the entropy can be computed from the (discretized) probability distribution,  $\{p_\alpha^i\}$ , for each feature  $i$ , as

$$H = \sum_{i,\alpha} -p_\alpha^i \log p_\alpha^i$$

5. **Select the SAE that minimizes the description length** i.e. the  $\varepsilon$ -MDL-optimal SAE.

## 320 F.2 Details for MNIST case study

321 For MNIST, we trained BatchTopK SAEs (Bussmann et al., 2024), typically for 1000+ epochs  
 322 until the test reconstruction loss converged or stopping early in cases of overfitting. Our desired  
 323 MSE tolerance was 0.0150. Discretizing the floats to roughly 5 bits per nonzero float gave an  
 324 average change in MSE of  $\approx 0.0001$ , which was roughly the scale over which MSE varied for the  
 325 hyperparameters used.

326 Gao et al. (2024) find that as the SAE width increases, there's a point where the number of dead  
 327 features starts to rise. In our experiments, we noticed that this point seems to be at a similar point to  
 328 where the description length starts to increase as well, although we did not test this systematically  
 329 and this property may be somewhat dataset dependent.

## 330 G Hierarchical features allow for more efficient coding schemes

331 Often features are semantically or causally related and this should allow for more efficient coding  
 332 schemes. For example, consider the hierarchical concepts "Animal" ( $A$ ) and "Bird" ( $B$ ). Since all  
 333 birds are animals, the "Animal" feature will always be active when the "Bird" feature is active. A  
 334 conventional SAE would represent these as separate feature vectors, one for "Bird" ( $B$ ) and one for  
 335 "Generic Animal" ( $A \wedge \neg B$ ), that are never active together, as shown in fig. 5. This setup has a low  
 336  $L_0$ , equal to the probability of "Animal",  $p_A$ , since something is a bird, a generic animal, or neither.

337 An alternative approach would be to define a variable length coding scheme (Salomon, 2007). For  
 338 example, one might consider first sending the activation for "Animal" ( $A$ ) and only if "Animal" is  
 339 active, sending the activation for "Animal is a Bird" ( $B|A$ ). Now the description length is given as  
 340  $DL = H(p_A) + p_A H(p_{B|A})$  which is always fewer bits compared to the conventional SAE with  
 341  $DL = H(p_A - p_B) + H(p_B)$ , (see the phase diagram in fig. 5). The overall  $L_0$  however is higher  
 342 because sometimes two activations are nonzero at the same time, so  $L_0 = p_A + p_{B|A}$ .

343 This case illustrates the potential to reduce description length by matching the SAE architecture  
 344 more closely to the hierarchical and causal structure of the data distribution. We also see another  
 345 case where optimising for sparsity differs to the MDL approach - hierarchical structures of the type  
 346 described above are never beneficial when optimising for sparsity but when thinking in terms of  
 347 Description Length, there are clear benefits to using the semantic structure of the data.

348 **H Description lengths for hierarchical features**

349 Independent additivity of feature explanations also implies that the description length of the set of  
 350 activations,  $\{z_i\}$ , is the sum of the lengths for each feature  $DL(\{z_i\}) = \sum_i DL(z_i)$ . If we know  
 351 the distribution of the activations,  $p_i(z)$ , then it is possible to send the activations using an average  
 352 description length equal to the distribution’s entropy,  $DL(z_i) = H(p_i) \equiv \sum_{z \in Z} -p_i(z) \log_2 p_i(z)$ .  
 353 For directed trees, the average description length of a child feature would be the conditional entropy,  
 354  $DL_{\text{child}}(z_i) = H(p_i | \text{parent active})$ , which accounts for the fact that  $DL = 0$  when the parent is not  
 355 active. This is one reason that directed tree-style SAEs can potentially have smaller descriptions than  
 356 conventional SAEs.

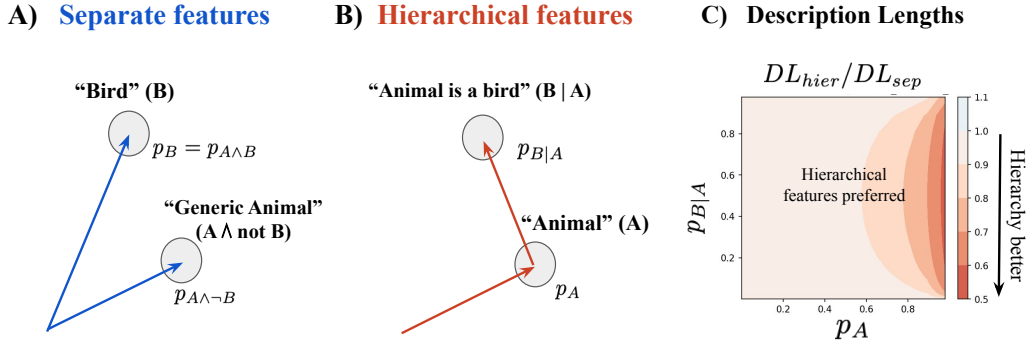


Figure 5: Two naturally hierarchical boolean features, such as "Animal" and "Bird", can be learned as separate mutually exclusive features (A) or in hierarchy (B) where the child feature can only be active if the parent feature is active, captured by the conditional probability  $p_{B|A}$ . C) The hierarchical case always has lower description length (DL) since the child feature’s activations need not be sent when the parent is not active.