
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 that SAE features representing significant line segments
12 are optimal, as opposed to SAEs with features for memorised digits from the
13 dataset or small digit fragments. We argue that using MDL rather than sparsity
14 may avoid potential pitfalls with naively maximising sparsity such as undesirable
15 feature splitting and that this framework naturally suggests new hierarchical SAE
16 architectures which provide more concise explanations.

17 1 Introduction

18 Sparse Autoencoders (SAEs) (Le, 2013; Makhzani and Frey, 2013) were developed to learn a
19 dictionary of sparsely activating features that describe a given dataset. They have recently become
20 popular tools for interpreting the internal activations of large foundation language models, often
21 finding human-understandable features (Sharkey et al., 2022; Huben et al., 2024; Bricken et al.,
22 2023b).

23 Interpretability, in particular human-understandability, is difficult to optimise for since ratings—from
24 humans or auto-interpretability methods (Bills et al., 2023)—are not differentiable at training time
25 and often cannot be efficiently obtained. Researchers often use sparsity, the number of nonzero
26 feature activations as measured by the L_0 norm, as a proxy for interpretability. SAEs are typically
27 trained with an additional L_1 penalty in their loss function to promote sparsity.

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

33 We introduce this information theoretic view by describing how SAEs can be used in a communication
34 protocol to transmit neural activations. We then argue that interpretability requires explanations
35 to have the property of independent additivity, which allows individual features to be interpreted

36 separately and discuss SAE architectures that are compatible with this property. We find that sparsity
37 (i.e. minimizing L_0) is a key component of minimizing description length but there are cases where
38 sparsity and description length diverge. In these cases, minimizing description length directly gives
39 more intuitive results. We demonstrate our approach empirically by finding the Minimal Description
40 Length solution for SAEs trained on the MNIST dataset.

41 2 SAEs are communicable explanations

42 SAEs aim to provide explanations of neural activations in terms of "features"¹. Here we reformulate
43 SAEs as solving a communication problem: suppose that we would like to transmit the neural
44 activations x to a friend with some tolerance ε , either in terms of the reconstruction error or change in
45 the downstream cross-entropy loss. Using the SAE as an encoding mechanism, we can approximate
46 the representation of the activations in two parts. *First*, we send them the SAE encodings of the
47 activations $z = Enc(x)$. *Second*, we send them a decoder network $Dec(\cdot)$ that recompiles these
48 activations back to (some close approximation of) the neural activations, $\hat{x} = Dec(z)$.

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

53 **Definition 2.1** *An explanation e of some phenomena p is a statement $e(p)$ for which knowing $e(p)$
54 gives some information about p . An explanation is typically a natural language statement².*

55 The description length (DL) of an explanation is the number of bits needed to transmit the explanation.
56 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
57 $O(1)$ in the dataset size so the first term dominates in the large data regime.

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

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

66 3 Interpretability requires independent additivity

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

70 SAE features can be interpreted either as **causal results** of the model inputs (which we can see
71 by analyzing feature activation patterns) or they can be interpreted as **causes** of the model outputs
72 (which we can see through conducting interventions on the features and seeing the downstream
73 effects). In both cases, we want to be able to understand each SAE feature independently, without
74 needing to control for the activations of the other features. If all the feature activations are causally
75 entangled—as is the case for the dense neural activations themselves—then they are not interpretable.
76 Note that for D features there are $O(D^2)$ pairs of features and $\sum_i^K \binom{D}{i}$ possible sets of features

¹Here 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.

²We 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.

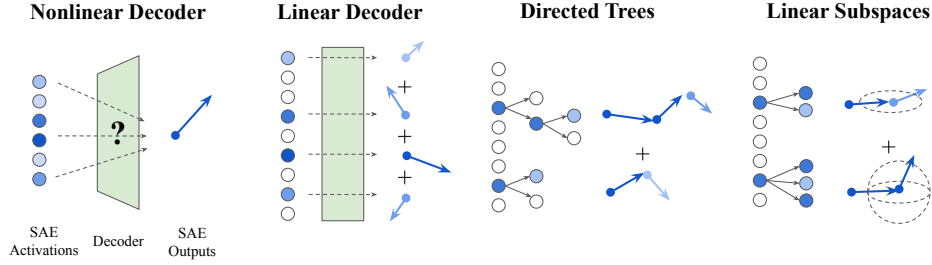


Figure 1: Examples of different SAE architectures. All but nonlinear decoders are compatible with independent additivity as feature activations correspond to adding a separate vector to the output. Architectures with directed tree decoders or which allow for vectors lying within a subspace are potentially more communication efficient since a child node can only be active if its parent node is active.

77 which is much too large for humans to hold in working memory. So for feature explanations to be
 78 human-understandable we cannot have the all the features being entangled such that understanding a
 79 single concept requires understanding arbitrary feature interactions.

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

83 **Definition 3.1** *Independent Additivity: An explanation e based on a vector of feature activations*
 84 $\vec{z} = \sum_i \vec{z}_i$ *is independently additive if $e(\vec{z}) \approx \sum_i e(\vec{z}_i)$. We say that a set of features z_i are*
 85 *independently additive if they can be understood independently of each other and the explanation of*
 86 *the sum of the features is the sum of the explanations of the features*³.

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

93 It is immediate from the definition that Independent Additivity holds for linear decoders however, we
 94 note that this condition also allows for more general decoder architectures. For example, features can
 95 be arranged to form a collection of directed trees, shown in fig. 1, where arrows represent the property
 96 "the child node can only be active if the parent node is active"⁴. Here each feature still corresponds to
 97 its own vector direction in the decoder. Since each child feature has a single path to its root feature,
 98 there are no interactions to disentangle and the independent additivity property still holds, in that
 99 each *tree* can be understood independently in a way that's natural for humans to understand, as a
 100 multi-dimensional feature. An advantage of the directed-tree SAE decoder structure is that it can be
 101 more description-length efficient as shown in fig. 5.

102 Independent additivity of feature explanations also implies that the description length of the set of
 103 activations, $\{z_i\}$, is the sum of the lengths for each feature $DL(\{z_i\}) = \sum_i DL(z_i)$. If we know
 104 the distribution of the activations, $p_i(z)$, then it is possible to send the activations using an average
 105 description length equal to the distribution's entropy, $DL(z_i) = H(p_i) \equiv \sum_{z \in Z} -p_i(z) \log_2 p_i(z)$.
 106 For directed trees, the average description length of a child feature would be the conditional entropy,
 107 $DL_{\text{child}}(z_i) = H(p_i | \text{parent active})$, which accounts for the fact that $DL = 0$ when the parent is not

³Note 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 ($\hat{x} = \text{Dec}(\vec{z}) = \sum_i \text{Dec}(z_i)$).

⁴In 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.

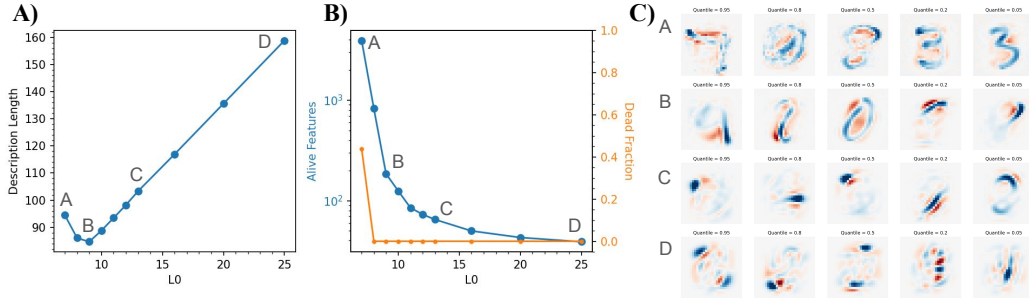


Figure 2: Finding the minimal description length (MDL) solution for SAEs trained on MNIST. A) Description length vs sparsity (L_0) for a set of hyperparameters with the same reconstruction error. B) Plot of the number of alive features as a function of sparsity (L_0). C) A random sample of SAE features at the 95th, 80th, 50th, 20th, and 5th percentiles of feature density respectively.

108 active. This is one reason that directed tree-style SAEs can potentially have smaller descriptions than
 109 conventional SAEs.

110 4 SAEs should be sparse, but not too sparse

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

115 The description length for a set of SAE activations (under independent additivity) with distribution
 116 $p(z)$ is given by $H(p) = \sum_{z \in Z} -p(z) \log_2 p(z)$. For exposition, consider a simpler formulation
 117 where we directly consider the bits needed without prior knowledge of the distributions. For a set of
 118 feature activations with L_0 nonzero elements out of D dictionary features, an upper bound on the
 119 description length is

$$DL \lesssim L_0(B + \log_2 D) \quad (1)$$

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

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

130 5 MDL-SAEs find interpretable and composable features for MNIST

131 Lee (2001) describe the classical method for using the Minimal Description Length (MDL) criteria
 132 for model selection. Here we choose between model hyperparameters (in particular the SAE width
 133 and expected L_0) for the optimal SAE. Our algorithm for finding the MDL-SAE solution and details
 134 for this case study are given in Appendix A.

135 We trained SAEs on the MNIST dataset of handwritten digits (LeCun et al., 1998) and find the set of
 136 hyperparameters resulting in the same test MSE. We see three basic regimes:

- 137 • **High L_0 , narrow SAE width** (C, D in fig. 2): Here, the description length (DL) is linear
 138 with L_0 , suggesting that the DL is dominated by the number of bits needed to represent the
 139 L_0 nonzero floats. The features appear as small sections of digits that could be relevant to
 140 many digits (C) or start to look like dense features that one might obtain from PCA (D).

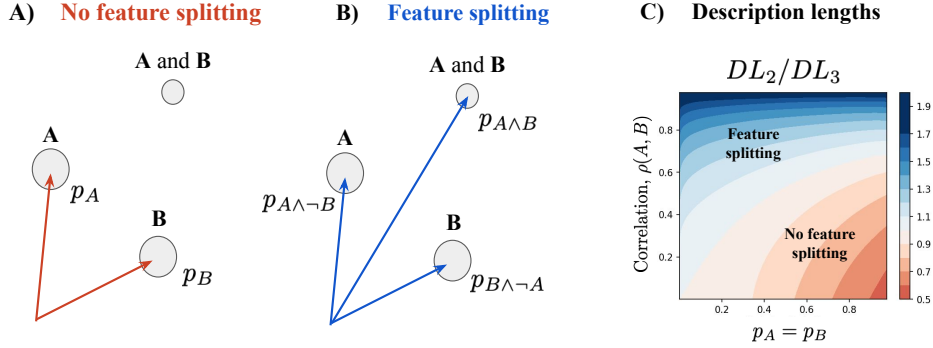


Figure 3: A toy model of undesirable feature splitting. The SAE can learn two boolean features without feature splitting (A) or three mutually exclusive boolean features with feature splitting (B) which always has lower L_0 . Minimizing description length provides a decision boundary (C) for when feature splitting is preferred or not.

- **Low L_0 , wide SAE width** (A in fig. 2): Though L_0 is small, the DL is large because as the SAE becomes wider, additional bits are required to specify which activations are nonzero. The features appear closer to being full digits, i.e. similar to samples from the dataset.
- **The MDL solution** (B in fig. 2): There’s a balance between the two contributions to the description length. The features appear like longer line segments or strokes for digits, but could apply to multiple digits.

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

6 Optimising for MDL can reduce undesirable feature splitting

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

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

No Feature Splitting: Say that the SAE only learns two boolean feature vectors, v_A and v_B , as shown in fig. 3. It is still capable of reconstructing $A \wedge B$ as the sum $v_A + v_B$. The L_0 would simply be the expectation of the boolean activations, so $L_0 = p_A + p_B$ and the description length would be $DL = H(p_A) + H(p_B)$ where $H(p)$ is the entropy of a Bernoulli variable with probability p .

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

Even though feature splitting always results in a lower L_0 , it does not always result in the smallest description length. The phase diagram in fig. 3 shows the case where $p_A = p_B$. If the correlation coefficient ρ between A and B is small then representing only A and B , but not $A \wedge B$, takes fewer bits so the preferred solution avoids feature splitting. However, if the correlation is large, then feature splitting is preferred since $A \wedge B$ occurs frequently enough that explicitly representing it reduces the description length. In this way, minimizing description length can limit the amount of undesirable

176 feature splitting and gives us a concrete decision criteria to understand when we might expect feature
177 splitting.

178 7 Hierarchical features allow for more efficient coding schemes

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

185 An alternative approach would be to define a variable length coding scheme (Salomon, 2007). For
186 example, one might consider first sending the activation for "Animal" (A) and only if "Animal" is
187 active, sending the activation for "Animal is a Bird" ($B|A$). Now the description length is given as
188 $DL = H(p_A) + p_A H(p_{B|A})$ which is always fewer bits compared to the conventional SAE with
189 $DL = H(p_A - p_B) + H(p_B)$, (see the phase diagram in fig. 5). The overall L_0 however is higher
190 because sometimes two activations are nonzero at the same time, so $L_0 = p_A + p_{B|A}$.

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

196 8 Related Work

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

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

207 9 Conclusion

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

215 Historically, evaluating SAEs for interpretability has been difficult without human interpretability
216 ratings studies, which can be labour intensive and expensive. We propose that operationalising inter-
217 pretability as efficient communication can help in creating principled evaluations for interpretability,
218 requiring less subjective and expensive SAE metrics. We would be excited about future work which
219 explores to what extent variants in SAE architectures can decrease the MDL of communicated latent
220 feature activations. In particular, we suggest that exploiting causal structure inherent in the data
221 distribution may be important to efficient coding.

222 References

- 223 S. Bills, N. Cammarata, D. Mossing, H. Tillman, L. Gao, G. Goh, I. Sutskever, J. Leike, J. Wu,
224 and W. Saunders. Language models can explain neurons in language models. [https://](https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html)
225 openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html, 2023.
- 226 J. Bloom. Open source sparse autoencoders for all residual stream layers of gpt2-small. *AI Align-*
227 *ment Forum*, 2024. URL [https://www.alignmentforum.org/posts/f9EgfLSurAiqRjySD/](https://www.alignmentforum.org/posts/f9EgfLSurAiqRjySD/open-source-sparse-autoencoders-for-all-residual-stream)
228 [open-source-sparse-autoencoders-for-all-residual-stream](https://www.alignmentforum.org/posts/f9EgfLSurAiqRjySD/open-source-sparse-autoencoders-for-all-residual-stream).
- 229 T. Bricken, J. Batson, A. Templeton, A. Jermyn, T. Henighan, and C. Olah. Features as the simplest
230 factorization. *Transformer Circuits Thread*, 2023a. URL [https://transformer-circuits.](https://transformer-circuits.pub/2023/may-update/index.html#simple-factorization)
231 [pub/2023/may-update/index.html#simple-factorization](https://transformer-circuits.pub/2023/may-update/index.html#simple-factorization).
- 232 T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner, C. Anil, C. Denison,
233 A. Askell, R. Lasenby, Y. Wu, S. Kravec, N. Schiefer, T. Maxwell, N. Joseph, Z. Hatfield-
234 Dodds, A. Tamkin, K. Nguyen, B. McLean, J. E. Burke, T. Hume, S. Carter, T. Henighan,
235 and C. Olah. Towards monosemanticity: Decomposing language models with dictionary learn-
236 ing. *Transformer Circuits Thread*, 2023b. URL [https://transformer-circuits.pub/2023/](https://transformer-circuits.pub/2023/monosemantic-features/index.html)
237 [monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 238 B. Bussmann, P. Leask, and N. Nanda. Batchtopk: A simple improvement for topk-saes. *AI Align-*
239 *ment Forum*, 2024. URL [https://www.alignmentforum.org/posts/Nkx6yWZNbAsfvic98/](https://www.alignmentforum.org/posts/Nkx6yWZNbAsfvic98/batchtopk-a-simple-improvement-for-topk-saes)
240 [batchtopk-a-simple-improvement-for-topk-saes](https://www.alignmentforum.org/posts/Nkx6yWZNbAsfvic98/batchtopk-a-simple-improvement-for-topk-saes).
- 241 L. Chan, R. Agrawal, A. Garriga-Alonso, and J. Gross. Compact proofs of model performance via
242 mechanistic interpretability. *AI Alignment Forum*, 2024. URL [https://www.alignmentforum.](https://www.alignmentforum.org/posts/bRsKimQcPTX3tNNJZ)
243 [org/posts/bRsKimQcPTX3tNNJZ](https://www.alignmentforum.org/posts/bRsKimQcPTX3tNNJZ).
- 244 L. Gao, T. D. la Tour, H. Tillman, G. Goh, R. Troll, A. Radford, I. Sutskever, J. Leike, and J. Wu.
245 Scaling and evaluating sparse autoencoders. *arXiv preprint arXiv:2406.04093*, 2024.
- 246 P. D. Grünwald. *The minimum description length principle*. MIT press, 2007.
- 247 R. Huben, H. Cunningham, L. R. Smith, A. Ewart, and L. Sharkey. Sparse autoencoders find highly
248 interpretable features in language models. In *The Twelfth International Conference on Learning*
249 *Representations*, 2024. URL <https://openreview.net/forum?id=F76bwRSLek>.
- 250 Q. V. Le. Building high-level features using large scale unsupervised learning. In *2013 IEEE*
251 *international conference on acoustics, speech and signal processing*, pages 8595–8598. IEEE,
252 2013.
- 253 Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document
254 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 255 T. C. Lee. An introduction to coding theory and the two-part minimum description length principle.
256 *International statistical review*, 69(2):169–183, 2001.
- 257 D. J. MacKay. *Information theory, inference and learning algorithms*. Cambridge university press,
258 2003.
- 259 A. Makhzani and B. Frey. K-sparse autoencoders. *arXiv preprint arXiv:1312.5663*, 2013.
- 260 C. Olah. Circuits updates - july 2024, linear representations. *Transformer Circuits Thread*, 2024.
261 <https://transformer-circuits.pub/2024/july-update/index.html#linear-representations>.
- 262 I. Ramirez and G. Sapiro. An mdl framework for sparse coding and dictionary learning. *IEEE*
263 *Transactions on Signal Processing*, 60(6):2913–2927, 2012.
- 264 D. Salomon. *Variable-length codes for data compression*. Springer Science & Business Media, 2007.
- 265 L. Sharkey, D. Braun, and B. Millidge. Taking features out of su-
266 perposition with sparse autoencoders. *AI Alignment Forum*, 2022.
267 URL [https://www.alignmentforum.org/posts/z6QQJbtpkEAX3Aojj/](https://www.alignmentforum.org/posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposition)
268 [interim-research-report-taking-features-out-of-superposition](https://www.alignmentforum.org/posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposition).

269 **A SAE communication protocol**

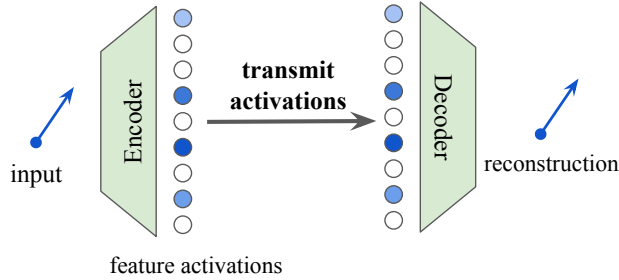


Figure 4: A schematic showing a sparse autoencoder (SAE) being used to communicate an input by transmitting the encoded activations and decoding them into a reconstruction of the input.

270 **B Comparison of GPT-2 SAE hyperparameters**

- 271 • **Reasonable SAEs:** Bloom (2024)’s open-source SAEs for GPT-2 layer 8 have $L_0 = 65$,
272 $D = 25,000$. Given $B = 7$ bits per nonzero float (8-bit quantization with the sign fixed to
273 positive), the description length per input token is 1405 bits.
- 274 • **Dense Activations:** A dense representation that still satisfies independent additivity would
275 be to send the neural activations directly instead of training an SAE. GPT-2 has a model size
276 of $d = 768$, the description length is simply $DL = B d = 5376$ bits per token.
- 277 • **One-hot encodings:** At the sparse extreme, our dictionary has a row for each neural
278 activation in the dataset, so $L_0 = 1$ and $D = (\text{vocab size})^{\text{seq len}}$. GPT-2 has a vocab size of
279 50,257 and the SAEs are trained 128 token sequences. All together this gives $DL = 13,993$
280 bits per token.

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

285 Note the difference here from choosing models based on the reconstruction loss vs sparsity (L_0)
286 Pareto frontier. When minimising L_0 , we are encouraging decreasing L_0 and increasing D until
287 $L_0 = 1$. Under the MDL model selection paradigm we are typically able to discount trivial solutions
288 like a one-hot encoding of the input activations and other extremely sparse solutions which make the
289 reconstruction algorithm analogous to a k-Nearest Neighbour classifier.

290 **C Details on determining the MDL-SAE**

291 **C.1 Algorithm**

- 292 1. **Specify a tolerance level, ϵ , for the loss function.** The tolerance ϵ is the maximum allowed
293 value for the loss, either the reconstruction loss (MSE for the SAE) or the model’s cross-
294 entropy loss when intervening on the model to swap in the SAE reconstructions in place of
295 the clean activations. For small datasets using a reconstruction, the test loss should be used.
- 296 2. **Train a set of SAEs within the loss tolerance.** It may be possible to simplify this task by
297 allowing the sparsity parameter to also be learned.
- 298 3. **Find the effective precision needed for floats.** The description length depends on the float
299 quantisation. We typically reduce the float precision until the change in loss results in the
300 reconstruction tolerance level is exceeded.

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

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

301 5. **Select the SAE that minimizes the description length** i.e. the ε -MDL-optimal SAE.

302 **C.2 Details for MNIST case study**

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

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

312 **D Description lengths for hierarchical features**

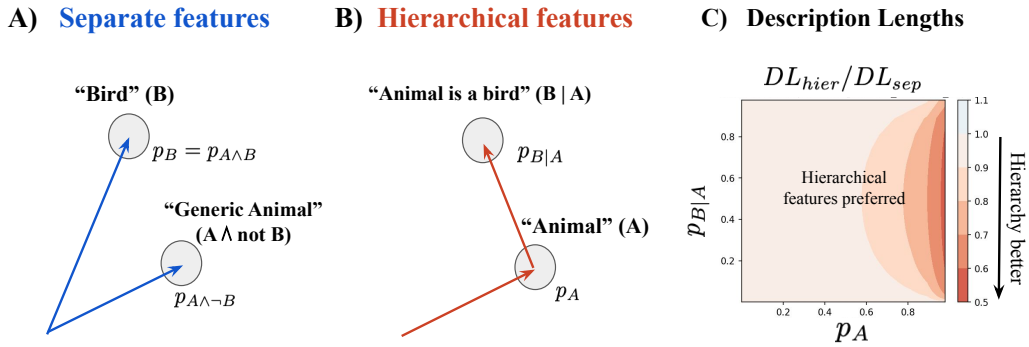


Figure 5: Two naturally hierarchical boolean features, such as "Animal" and "Bird", can be learned as separate mutually exclusive features (A) or in hierarchy (B) where the child feature can only be active if the parent feature is active, captured by the conditional probability $p_{B|A}$. C) The hierarchical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