
Priors in Time: A Generative View of Sparse Autoencoders for Sequential Representations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Sparse Autoencoders (SAEs) are widely used to decompose neural network representations into interpretable concepts. Despite their success, SAEs often fail to
2 capture all relevant concepts, raising the question of what assumptions underlie
3 their limitations. We show that these challenges arise from a mismatch between
4 the true data distribution and the implicit priors encoded in SAE architectures
5 and sparsity regularizers. Taking language model representations as a case study,
6 we demonstrate that these activations exhibit rich temporal structure—such as
7 systematic growth in concept dimensionality, context-dependent correlations, and
8 non-stationarity over time—that conflicts with SAE priors. Through experiments,
9 we highlight how this mismatch leads to characteristic SAE pathologies, including
10 degraded concept recovery and reconstruction quality over time. Our results point
11 toward the need for new SAE designs that incorporate inductive biases aligned
12 with temporal dynamics in sequential data.
13

14 1 Introduction

15 Scaling of neural networks [1–4] has been argued to enable
16 better approximations of the data distribution [5–8], leading
17 to learning of cognitive abilities generally associated with humans [9–12].
18 Since computations underlying these capabilities occur via intermediate representations, one can intuitively
19 expect the representations encode concepts relevant to such abilities (i.e., latent factors of the generative process) within
20 them [13–19]. With the dual goal of better understanding human cognition and controlling models’, this argument has motivated
21 recent interpretability work on identifying what concepts underlie a model’s abilities. Specifically, such works build on
22 hypothesized computational models of how concepts are encoded in neural network representations, motivating tools for
23 unsupervised extraction of a dictionary of vectors that (ideally) correspond to meaningful concepts. The prime example of such
24 computational models and corresponding tools includes the linear representation hypothesis (LRH) [20–23] and sparse autoencoders (SAEs) [24–27].

25
26
27
28
29
30
31
32 **This work.** If neural network representations indeed encode concepts underlying the distribution, any time a computational model is proposed for how concepts
33 organize in representations, one ends up making implicit assumptions about the data distribution. For example, motivated by literature
34 on representation steering, SAEs assume model behavior can be manipulated along a given concept without interfering with other ones. This amounts to assuming
35 statistical independence of concepts
36

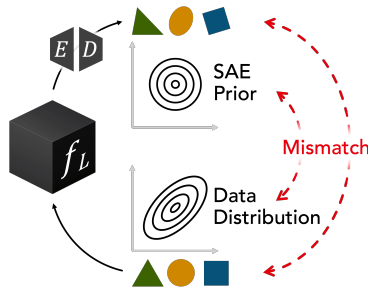


Figure 1: **SAEs assume a prior on concepts that generate data.** When this prior deviates from true data structure, SAEs end up learning different/distorted concepts.

involved in the generative process. For natural domains, e.g., language, such an assumption can be restrictively strict, since even simple linguistic concepts such as tense and prepositions will be correlated. The mismatch between such implicit assumptions versus the true generative process can thus constrain what concepts an SAE is able to identify from representations (Fig. 1). We aim to formalize this problem: specifically, we make explicit the assumptions about the data-generating process implicit in SAEs and demonstrate the consequences thereof in the specific setting of interpreting language models. In brief, our contributions can be summarized as follows.

- **Characterizing prior assumptions of existing SAEs.** By adopting a Bayesian perspective on SAE training objectives, we prove that SAEs implicitly assume a prior distribution over concepts.
- **Revealing temporal structure in LLM representations.** Contextualizing our theory with respect to large language models (LLMs), we demonstrate rich temporal dynamics in an LLM’s representations—including systematic growth in the number of active concepts, context-sensitive correlations, and non-stationarity—that are inconsistent with implicit assumptions made by SAEs.
- **Demonstrating SAE limitations under prior-data mismatch.** The mismatch between prior assumptions of SAEs and actual structure of LLM representations suggests SAEs will exhibit poor reconstruction error with increasing sequence length. We empirically verify this claim with off-the-shelf, pretrained SAEs.

Overall, our results suggest interpretability should be driven by the behavior one is trying to explain—else, mismatch in structure of the distribution describing the behavior and tools used can lead to pathological results [28–32].

2 Priors of Sparse Autoencoders

Let bold, lowercase letters represent vectors (e.g., \mathbf{z}). Subscripts on vectors denote different samples (e.g., \mathbf{z}_i), while superscripts denote the index within the vector, leading to a scalar (e.g., z^k). We denote model activations by $\mathbf{x} \in \mathbb{R}^n$, SAE latents (sparse code) by $\mathbf{z} \in \mathbb{R}^M$, and the dictionary by $\mathbf{D} \in \mathbb{R}^{n \times M}$ (M is the dictionary size). SAEs solve the following optimization problem, which is equivalent to constrained dictionary learning ([29]): $\arg \min_{\mathbf{D}, \mathbf{z}} \sum_{i=1}^N \frac{1}{N} (\|\mathbf{x}_i - \mathbf{D}\mathbf{z}_i\|_2^2 + \lambda \mathcal{R}(\mathbf{z}_i))$, where $\forall k, \mathbf{z}_k = f_{\text{SAE}}(\mathbf{x}_k)$. This formulation affords a Bayesian lens on SAEs as a MAP (maximum-a-posteriori) estimation problem for the sparse codes $\{\mathbf{z}_i\}$, given data $\{\mathbf{x}_i\}$, where the first term (MSE) captures the log likelihood, while the latter (sparsity penalty) is the log prior. This perspective leads us to the following theorem about priors of various SAEs.

Theorem 2.1 (SAE Priors on Sparse Code). *Let $S_t = \text{supp}(\mathbf{z}_t) = \{k : z_t^k > 0\}$ be the set of active latents in the sparse code \mathbf{z} at time t , and $n_t = |S_t|$ be the cardinality of S_t (the number of active latents). Each SAE imposes a prior distribution on the sparse code \mathbf{z} , arising from its sparsity penalty $\mathcal{R}(\mathbf{z})$ or implicit conditions imposed on the sparse code. These conditions are highlighted in Table 1.*

Table 1: Priors over concept interactions and dynamics for various SAEs

$f_{\text{SAE}}, \mathcal{R}(\mathbf{z})$	Across-Concept Prior (interaction)	Across-time Prior (dynamics)
ReLU, L_1 -norm	$z_t^1, \dots, z_t^M \stackrel{\text{i.i.d.}}{\sim} \text{Laplace}(0, \cdot)$	$\mathbf{z}_1, \dots, \mathbf{z}_t \stackrel{\text{i.i.d.}}{\sim} P^{(1)}$
TopK	$z_t^{i_1}, \dots, z_t^{i_{n_t}} \mid S_t \stackrel{\text{i.i.d.}}{\sim} U(0, \cdot) \forall i. \in S_t,$ $S_t \sim U([M]^K)$	$(\mathbf{z}_1, S_1), \dots, (\mathbf{z}_t, S_t) \stackrel{\text{i.i.d.}}{\sim} P^{(2)},$ $S_1, \dots, S_t \stackrel{\text{i.i.d.}}{\sim} U([M]^K)$
JumpReLU, L_0 -norm	$z_t^{i_1}, \dots, z_t^{i_{n_t}} \mid S_t \stackrel{\text{i.i.d.}}{\sim} U(0, \cdot) \forall i. \in S_t,$ $S_t \mid n_t \sim U([M]^{n_t})$	$(\mathbf{z}_1, S_1, n_1), \dots, (\mathbf{z}_t, S_t, n_t) \stackrel{\text{i.i.d.}}{\sim} P^{(3)},$ $n_1, \dots, n_t \stackrel{\text{i.i.d.}}{\sim} P^{(4)}$
BatchTopK	$z_t^{i_1}, \dots, z_t^{i_{n_t}} \mid S_t \stackrel{\text{i.i.d.}}{\sim} U(0, \cdot) \forall i. \in S_t,$ $S_t \mid n_t \sim U([M]^{n_t})$	$(\mathbf{z}_1, S_1, n_1), \dots, (\mathbf{z}_t, S_t, n_t) \stackrel{\text{i.i.d.}}{\sim} P^{(5)},$ $n_1, \dots, n_t \stackrel{\text{i.i.d.}}{\sim} P^{(6)}, \mathbb{E}[n_t] = K$

71

Implications for concepts: ReLU SAE assumes an independence prior over concepts at each time t , as well as independence of concepts over time. TopK, BatchTopK and JumpReLU SAEs assume conditional independence of concepts conditioned on the set of active concepts at each time, and independence of concepts along with their sparsity over time.

(a) There stood a dark-haired boy, his Gryffindor scarf blazing red and gold against the night, round glasses framing the scar that marked him forever, wand in hand – this was Harry Potter.

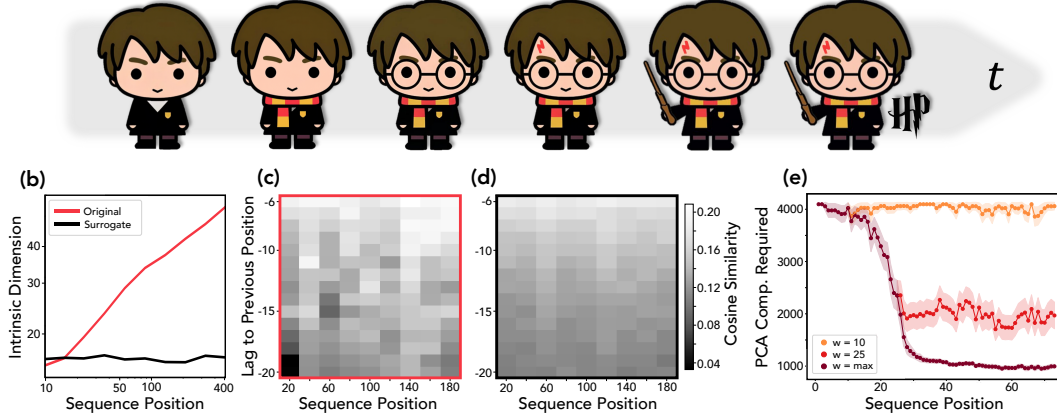


Figure 2: **Language model representations have rich temporal structure:** Illustration (in (a)) of and evidence (b–e) for temporal properties of LLM representations. (a) An illustrative, toy example of Harry Potter demonstrating the evolution of character attributes over a sentence. (b) Intrinsic dimension (measured using a U-statistic of cosine similarities between vectors) as a function of time. (c) Autocorrelation function of representations as a function of time (x-axis) and lag/ shift (y-axis); (d) baseline autocorrelation for a stationary process. (e) Number of PCA components (basis constructed from past context) needed to explain 90% variance in representations as a function of time.

76 3 Temporal Structure in LLM Representations

77 Natural language has a temporal structure, as illustrated in Fig. 2(a). The individual characters
 78 depict the reader’s growing knowledge across the given sentence about Harry Potter. There are three
 79 notable observations: (1) the number of concepts (attributes of the character) increases with time,
 80 with changes being incremental at each timestep; (2) the sum total of concepts at a given time t is
 81 strongly correlated with the past context, while only a few additions are novel; and (3) the sentence is
 82 non-stationary. Since a different concept is added at each step, the correlations between concepts at a
 83 given position and past positions are time-dependent.

84 As LLMs model natural language, we expect LLM representations to reflect such temporal structure.
 85 To this end, we investigate the temporal structure in LLM representations when processing webtext,
 86 to enable a comparison with the priors of existing SAEs discussed in Section 2. Specifically, we
 87 investigate the Llama-3.1-8B model representations on the Pile dataset [33]. Fig. 2(b)–(e) provide
 88 evidence for the non-stationarity of LLM representations. Panel (b) shows a steady increase in the
 89 intrinsic dimension of representations over time. This trend indicates an increase in the number of
 90 concepts, assuming concept representations are quasi-orthogonal [34, 30]. The autocorrelation in
 91 Panels (c) and (d) show the similarity between a representation at position t with earlier representa-
 92 tions. The autocorrelation pattern of LLM representations (c) is different across time t , indicating
 93 that representations are non-stationary. For comparison, Panel (d) shows that the autocorrelation of
 94 a stationary process is approximately uniform across the sequence. Appendix B provides further
 95 details on stationarity measures and the surrogate stationary process. Finally, in Panel (e), we project
 96 representations at a given time onto the PCA basis obtained using a context window with varying
 97 size. We note that 90% of the variance in representation can be explained using a context window of
 98 the past 50 tokens, strongly suggesting correlations between a representation and its context.

99 4 SAE Pathologies on Temporal Dynamics

100 The mismatch between SAE priors (Theorem 2.1) and the actual temporal structure of LLM represen-
 101 tations (Section 3) leads to three systematic failures, illustrated in Fig. 3.

102 **Sparsity dynamics** (Column a): Neither ReLU nor TopK SAEs capture the systematic growth in
 103 concept count observed in LLM representations (Fig. 2(b)). ReLU SAEs maintain nearly-constant
 104 sparsity throughout most of the sequence, with L_0 norm increasing only at late positions. TopK

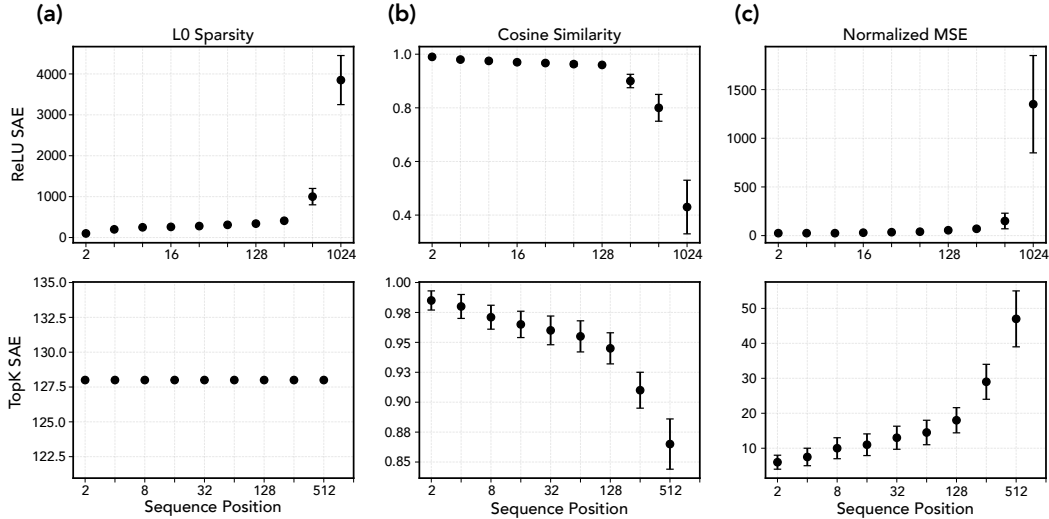


Figure 3: **SAEs fail to capture temporal structure in LLM representations:** Results from ReLU, TopK SAEs trained on GPT-2 model activations. (a) The sparsity of representations (L_0 (pseudo)-norm) as a function of time. (b) Cosine similarity between learnt features and true features decreases with time, making concept extraction progressively worse with time. (c) Normalized MSE increases with time, indicating worse reconstruction for later parts of a story.

105 SAEs, constrained by their architecture to exactly K active features, cannot adapt to varying sparsity
 106 levels at all.

107 **Concept recovery** (Column b): Cosine similarity between learned and true concept directions
 108 degrades progressively along the sequence for both SAE types, indicating deteriorating concept
 109 recovery—whereas an ideal SAE would maintain high similarity throughout.

110 **Reconstruction quality** (Column c): Normalized mean-squared error (NMSE) increases mono-
 111 tonically with sequence position, demonstrating that both SAE variants struggle increasingly to
 112 reconstruct representations with increasing sequence length.

113 5 Discussion

114 Our findings highlight fundamental limitations of existing SAEs when applied to sequential data.
 115 Most SAEs implicitly assume independence of concepts across time, yet our analysis shows that LLM
 116 representations exhibit rich temporal dynamics. These include a systematic growth in the number
 117 of active concepts, strong correlations with recent context, and non-stationarity in their temporal
 118 dependencies. This prior-data mismatch leads to systematic pathologies: in particular, we show
 119 that ReLU and TopK SAEs fail to capture the progressive increase in conceptual complexity and
 120 become increasingly ineffective at both concept recovery and representation reconstruction over
 121 time. We hypothesize that abstract, context-dependent phenomena—such as humor, sarcasm, or
 122 other higher-level semantic constructs—are especially unlikely to be recovered under current SAE
 123 formulations, as they emerge later in sequences and rely heavily on contextual correlations.

124 These results suggest the need for SAE variants with stronger inductive biases tailored to sequential
 125 representations. Future architectures may explicitly account for temporal correlations and incremental
 126 novelty in activations, for example by incorporating recurrence or self-attention over past concepts,
 127 and imposing sparsity penalty on incremental novel concepts rather than all concepts. Such modifi-
 128 cations could enable SAEs to better align with the true dynamics of LLM representations, thereby
 129 improving their ability to extract meaningful and temporally coherent concepts.

References

- 130
- 131 [1] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,
132 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
133 models. *arXiv preprint arXiv:2001.08361*, 2020.
- 134 [2] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
135 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.
136 Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- 137 [3] Jonathan S Rosenfeld. Scaling laws for deep learning. *arXiv preprint arXiv:2108.07686*, 2021.
- 138 [4] Corinna Cortes, Lawrence D Jackel, Sara Solla, Vladimir Vapnik, and John Denker. Learning
139 curves: Asymptotic values and rate of convergence. *Advances in neural information processing*
140 *systems*, 6, 1993.
- 141 [5] Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining
142 neural scaling laws. *Proceedings of the National Academy of Sciences*, 121(27):e2311878121,
143 2024.
- 144 [6] Alexander Maloney, Daniel A Roberts, and James Sully. A solvable model of neural scaling
145 laws. *arXiv preprint arXiv:2210.16859*, 2022.
- 146 [7] Blake Bordelon, Alexander Atanasov, and Cengiz Pehlevan. A dynamical model of neural
147 scaling laws. *arXiv preprint arXiv:2402.01092*, 2024.
- 148 [8] Daniel Wurgaft, Ekdeep Singh Lubana, Core Francisco Park, Hidenori Tanaka, Gautam Reddy,
149 and Noah D Goodman. In-context learning strategies emerge rationally. *arXiv preprint*
150 *arXiv:2506.17859*, 2025.
- 151 [9] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece
152 Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general
153 intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 154 [10] Google Deepmind. Advanced version of gemini with deep think
155 officially achieves gold-medal standard at the international math-
156 ematical olympiad. [https://deepmind.google/discover/blog/
157 advanced-version-of-gemini-with-deep-think-officially-achieves-gold-medal-standard-at-the-i-
158 2025.](https://deepmind.google/discover/blog/advanced-version-of-gemini-with-deep-think-officially-achieves-gold-medal-standard-at-the-i-)
- 159 [11] James WA Strachan, Dalila Albergo, Giulia Borghini, Oriana Pansardi, Eugenio Scaliti, Saurabh
160 Gupta, Krati Saxena, Alessandro Rufo, Stefano Panzeri, Guido Manzi, et al. Testing theory of
161 mind in large language models and humans. *Nature Human Behaviour*, 8(7):1285–1295, 2024.
- 162 [12] Laura Ruis, Akbir Khan, Stella Biderman, Sara Hooker, Tim Rocktäschel, and Edward Grefen-
163 stette. The goldilocks of pragmatic understanding: Fine-tuning strategy matters for implicature
164 resolution by llms. *Advances in Neural Information Processing Systems*, 36:20827–20905,
165 2023.
- 166 [13] Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic
167 development in deep neural networks. *Proceedings of the National Academy of Sciences*,
168 116(23):11537–11546, 2019.
- 169 [14] Joshua Engels, Eric J Michaud, Isaac Liao, Wes Gurnee, and Max Tegmark. Not all language
170 model features are one-dimensionally linear. *arXiv preprint arXiv:2405.14860*, 2024.
- 171 [15] Core Francisco Park, Andrew Lee, Ekdeep Singh Lubana, Yongyi Yang, Maya Okawa, Kento
172 Nishi, Martin Wattenberg, and Hidenori Tanaka. Iclr: In-context learning of representations,
173 2025.
- 174 [16] Bo Zhao, Maya Okawa, Eric J Bigelow, Rose Yu, Tomer Ullman, Ekdeep Singh Lubana, and
175 Hidenori Tanaka. Emergence of hierarchical emotion organization in large language models.
176 *arXiv preprint arXiv:2507.10599*, 2025.

- 177 [17] Michael Pearce, Elana Simon, Michael Byun, and Daniel Balsam. Finding the tree of life in evo
178 2. *Goodfire Research*, August 2025. Correspondence to michael@goodfire.ai.
- 179 [18] Kiho Park, Yo Joong Choe, Yibo Jiang, and Victor Veitch. The geometry of categorical and
180 hierarchical concepts in large language models. *arXiv preprint arXiv:2406.01506*, 2024.
- 181 [19] Alexander Modell, Patrick Rubin-Delanchy, and Nick Whiteley. The origins of representation
182 manifolds in large language models. *arXiv preprint arXiv:2505.18235*, 2025.
- [20] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna
Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse,
Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah.
Toy models of superposition. *Transformer Circuits Thread*, 2022. [https://transformer-circuits-pub/2022/toy_model/index.html](https://transformer-circuits.pub/2022/toy_model/index.html).
- [21] Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. Linear algebraic structure
184 of word senses, with applications to polysemy. *Transactions of the Association for Computational
185 Linguistics*, 6:483–495, 2018.
- [22] Jack Merullo, Noah A Smith, Sarah Wiegrefe, and Yanai Elazar. On linear representations and
187 pretraining data frequency in language models. *arXiv preprint arXiv:2504.12459*, 2025.
- [23] Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas,
189 Yonatan Belinkov, and David Bau. Linearity of relation decoding in transformer language models.
190 *arXiv preprint arXiv:2308.09124*, 2023.
- [24] Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya
192 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint
193 arXiv:2406.04093*, 2024.
- [25] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-
195 coders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*,
196 2023.
- [26] Bart Bussmann, Patrick Leask, and Neel Nanda. Batchtopk sparse autoencoders, 2024.
- [27] Senthoran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János
199 Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse
200 autoencoders, 2024.
- [28] David Chanin, James Wilken-Smith, Tomas Dulka, Hardik Bhatnagar, Satvik Golechha, and Joseph
202 Bloom. A is for absorption: Studying feature splitting and absorption in sparse autoencoders, 2025.
- [29] Sai Sumedh R Hindupur, Ekdeep Singh Lubana, Thomas Fel, and Demba Ba. Projecting assumptions:
204 The duality between sparse autoencoders and concept geometry. *arXiv preprint arXiv:2503.01822*,
205 2025.
- [30] Valérie Costa, Thomas Fel, Ekdeep Singh Lubana, Bahareh Tolooshams, and Demba Ba. From
207 flat to hierarchical: Extracting sparse representations with matching pursuit. *arXiv preprint
208 arXiv:2506.03093*, 2025.
- [31] Bart Bussmann, Noa Nabeshima, Adam Karvonen, and Neel Nanda. Learning multi-level features
210 with matryoshka sparse autoencoders. *arXiv preprint arXiv:2503.17547*, 2025.
- [32] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick
212 Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec,
213 Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen,
214 Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah.
215 Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer
216 Circuits Thread*, 2023. [https://transformer-circuits-pub/2023/monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- [33] Monology. The pile: Uncopyrighted subset. [https://huggingface.co/datasets/monology/
218 pile-uncopyrighted](https://huggingface.co/datasets/monology/pile-uncopyrighted), 2021. Based on the original Pile dataset by Gao et al.

- [24] Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. *arXiv preprint arXiv:2311.03658*, 2023.
- [25] Alireza Makhzani and Brendan Frey. K-sparse autoencoders. *arXiv preprint arXiv:1312.5663*, 2013.

222 **A Further Theory Results**

223 **A.1 Priors on the Sparse Code for various SAEs**

224 We first (re)state the loss function used in training SAEs:

$$\begin{aligned} & \arg \min_{\mathbf{D}, \mathbf{z}} \sum_{i=1}^N \frac{1}{N} (\|\mathbf{x}_i - \mathbf{D}\mathbf{z}_i\|_2^2 + \lambda \mathcal{R}(\mathbf{z}_i)) \\ & \text{s.t. } \forall k, \mathbf{z}_k = f_{SAE}(\mathbf{x}_k) \end{aligned}$$

225 As described in Section ??, the above loss can be viewed as the negative log posterior in a MAP
226 estimation problem for the sparse code $\{\mathbf{z}^{(i)}\}_i$ given data $\{\mathbf{x}^{(i)}\}_i$. In this perspective, the prior is
227 given as:

$$\log P(\mathbf{z}) \propto - \sum_{i=1}^N \mathcal{R}(\mathbf{z}_i) \quad (1)$$

228 **A.1.1 ReLU SAE**

229 The vanilla ReLU SAE ([32], [25]) is trained with the L_1 -norm penalty:

$$\mathcal{R}(\mathbf{z}) = \|\mathbf{z}\|_1. \quad (2)$$

230 The prior over \mathbf{z} for the above case is:

$$\log P(\mathbf{z}_1, \dots, \mathbf{z}_N) \propto - \sum_{i=1}^N \sum_{k=1}^M |z_i^k|, \quad (3)$$

$$\implies P(\mathbf{z}_1, \dots, \mathbf{z}_N) \propto \prod_{i=1}^N \left(\prod_{k=1}^M \exp -\nu |z_i^k| \right). \quad (4)$$

231 This joint distribution implies that for each sample i , different indices k are sampled i.i.d. from the
232 same distribution:

$$z_i^1, \dots, z_i^M \stackrel{\text{i.i.d.}}{\sim} \text{Laplace}(0, 1/\nu), \quad (5)$$

233 and different samples are all independently sampled from the same product Laplace distribution:

$$\mathbf{z}_1, \dots, \mathbf{z}_N \stackrel{\text{i.i.d.}}{\sim} \text{Laplace}^M(0, 1/\nu). \quad (6)$$

234 This concludes the proof for priors of ReLU SAE trained with L_1 norm sparsity penalty. \square

235 **A.1.2 TopK SAE**

236 The TopK SAE ([35], [24]) directly controls the sparsity of the representation \mathbf{z} by fixing it at
237 $\|\mathbf{z}\|_0 = K$, instead of imposing an explicit sparsity penalty $\mathcal{R}(\mathbf{z})$ in the loss function. The objective
238 function for TopK SAE is:

$$\arg \min_{\mathbf{D}, \mathbf{z}} \sum_{i=1}^N \frac{1}{N} \|\mathbf{x}_i - \mathbf{D}\mathbf{z}_i\|_2^2, \quad (7)$$

$$\text{s.t. } \forall j, \mathbf{z}_j = f_{TopK}(\mathbf{x}_j), \|\mathbf{z}_j\|_0 = K. \quad (8)$$

239 Since the fixed sparsity is a hard constraint that depends on \mathbf{z} alone (and not the data \mathbf{x}), we can use
240 Lagrange multipliers to reformulate it as an effective prior:

$$\arg \min_{\mathbf{D}, \mathbf{z}, \{\lambda_i\}} \sum_{i=1}^N \frac{1}{N} \left(\|\mathbf{x}_i - \mathbf{D}\mathbf{z}_i\|_2^2 + \lambda_i (\|\mathbf{z}_i\|_0 - K) \right), \quad (9)$$

$$\text{s.t. } \forall j, \mathbf{z}_j = f_{TopK}(\mathbf{x}_j). \quad (10)$$

241 The prior over \mathbf{z} for the above (effective) regularizer is:

$$\log P(\mathbf{z}_1, \dots, \mathbf{z}_N) \propto - \sum_{i=1}^N \lambda_i (|\|\mathbf{z}_i\|_0 - K|) \quad (11)$$

$$\implies P(\mathbf{z}_1, \dots, \mathbf{z}_N) \propto \prod_{i=1}^N \exp(-\lambda_i |\|\mathbf{z}_i\|_0 - K|) \quad (12)$$

242 Note that the above prior is finite for finite values of λ_i , but the overall objective optimizes over $\{\lambda_i\}$,
 243 resulting in a *hard* prior peaked at $\|\mathbf{z}_i\|_0 = K$ for each sample i .

244 The factorization over samples i implies mutual independence of $\mathbf{z}_1, \dots, \mathbf{z}_n$: $P(\mathbf{z}_1, \dots, \mathbf{z}_n) =$
 245 $\prod_{i=1}^N P(\mathbf{z}_i)$.

246 As defined in Theorem 2.1 (and restated here for convenience), let $S_i = \text{supp}(\mathbf{z}_i) = \{k : z_i^k >$
 247 $0\}$, $n_i = |S_i| = \|\mathbf{z}_i\|_0$ denote the active indices and their number (sparsity) respectively.

248 For individual samples \mathbf{z}_i , if we condition on the set of active indices S_i , the sparsity gets fixed since
 249 $\|\mathbf{z}_i\|_0 = |S_i| = n_i$, and the distribution becomes constant:

$$P(\mathbf{z}_i | S_i) = C \quad (13)$$

$$\implies z_i^\mu | S_i \sim \begin{cases} U(0, \kappa) & \mu \in S_i, \text{ and} \\ \delta_0 & \mu \notin S_i \end{cases} \quad (14)$$

$$z_i^{\mu_1}, \dots, z_i^{\mu_{|S_i|}} | S_i \stackrel{\text{i.i.d.}}{\sim} U(0, \kappa) \text{ for } \mu. \in S_i \quad (15)$$

250 where C, κ are appropriate constants.

251 Since $\{\mathbf{z}_i\}_i$ s are mutually independent, any measurable function of each is also independent. The
 252 indices of nonzero entries of \mathbf{z}_j , i.e., S_j is a measurable function since it is a map $S : \mathbb{R}_+^M \rightarrow 2^M$
 253 which is discrete valued, and pre images of each value—a set of nonzero indices—are measurable
 254 since they equal the cartesian products of the measurable sets $\{z = 0\}, \{z > 0\}$ over all indices.
 255 Hence, S_1, \dots, S_n are also independent.

256 Since $S_i = g(\mathbf{z}_i)$ and the distribution of \mathbf{z}_i depends only on $n_i = \|\mathbf{z}_i\|_0$ (Eq.11), the distribution of
 257 S_i will also depend only on n_i , becoming uniform when conditioned on n_i . In TopK SAE, $n_i = K$
 258 is a constant. Therefore, each $S_i \sim U([M]^K)$, and together with independence argued above,

$$S_1, \dots, S_N \stackrel{\text{i.i.d.}}{\sim} U([M]^K) \quad (16)$$

259 This completes the proof for the priors of TopK SAE. \square

260 A.1.3 BatchTopK SAE

261 BatchTopK SAE ([26]) is a modification of the TopK SAE. Instead of fixing sparsity like TopK,
 262 BatchTopK allows variable sparsity per input while fixing the mean sparsity over a batch at K . The
 263 objective function for BatchTopK SAE can equivalently be written as:

$$\arg \min_{\mathbf{D}, \mathbf{z}} \sum_{i=1}^N \frac{1}{N} \|\mathbf{x}_i - \mathbf{D}\mathbf{z}_i\|_2^2, \quad (17)$$

$$\text{s.t. } \forall j, \mathbf{z}_j = f_{\text{TopK}}(\mathbf{x}_j), \frac{1}{N} \sum_{j=1}^N \|\mathbf{z}_j\|_0 = K. \quad (18)$$

264 While BatchTopK imposes a mean sparsity per batch, for simplicity, we use the batch size to match the
 265 size of the entire dataset (WLOG). Smaller batch sizes can easily be incorporated by adding separate
 266 constraints, each over the entire batch (only leads to a change in constants—lagrange multipliers—in
 267 the analysis).

268 Following similar analysis as for TopK SAE (App. A.1.2), we can derive an equivalent prior over \mathbf{z}
 269 for BatchTopK SAE:

$$P(\mathbf{z}_1, \dots, \mathbf{z}_N) \propto \prod_{i=1}^N \exp(-\lambda |\|\mathbf{z}_i\|_0 - K|) \quad (19)$$

270 The sparse codes for different samples $\{z_i\}_i$ are thus sampled i.i.d. from a distribution that only
 271 depends on the sparsity penalty. While this prior looks very similar to the prior of TopK SAE, the
 272 difference is that in TopK, the fixed sparsity constraint is imposed per sample, leading to a different
 273 Lagrange multiplier λ_i per sample to optimize over, while in BatchTopK, we have a common
 274 multiplier λ over all examples in a batch (with multiple batches, we will have one multiplier per
 275 batch), which is then optimized over to ensure that average sparsity per batch constraint is met.

276 Similar to the analysis for the TopK SAE, we get the following prior over different latents per sample:

$$z_i^\mu | S_i \sim \begin{cases} U(0, \kappa) & \mu \in S_i \\ \delta_0 & \mu \notin S_i \end{cases}, \text{ and} \quad (20)$$

$$z_i^{\mu_1}, \dots, z_i^{\mu_{|S_i|}} | S_i \stackrel{\text{i.i.d.}}{\sim} U(0, \kappa) \text{ for } \mu. \in S_i \quad (21)$$

277 The active indices S_i are sampled uniformly conditioned on the number of active indices n_i :

$$S_i | n_i \sim U([M]^{n_i}) \quad (22)$$

278 The number of active latents n_i are themselves sampled i.i.d. (since $n_i = \tilde{g}(z_i)$ and $\{z_i\}_i$ are i.i.d.)
 279 from a distribution whose mean is fixed:

$$n_1, \dots, n_N \stackrel{\text{i.i.d.}}{\sim} P, \text{ s.t. } \mathbb{E}[n.] = K \quad (23)$$

280 This completes the derivation for the BatchTopK prior. \square

281 **A.1.4 JumpReLU SAE**

282 JumpReLU SAE ([27]) is trained with the L_0 (pseudo-)norm regularizer. This leads to the following
 283 optimization problem:

$$\begin{aligned} \arg \min_{D, z} \sum_{i=1}^N \frac{1}{N} (\|x_i - Dz_i\|_2^2 + \lambda \|z_i\|_0) \\ \text{s.t. } \forall k, z_k = f_{\text{JumpReLU}}(x_k) \end{aligned}$$

284 This objective is equivalent to the following prior over z :

$$P(z_1, \dots, z_N) \propto \prod_{i=1}^N \exp(-\eta \|z_i\|_0) \quad (24)$$

285 Noting the similarity with the TopK/ BatchTopK cases, we use the same analysis to derive the
 286 following conditions:

$$z_i^\mu | S_i \sim \begin{cases} U(0, \kappa) & \mu \in S_i \\ \delta_0 & \mu \notin S_i \end{cases}, \text{ and} \quad (25)$$

$$z_i^{\mu_1}, \dots, z_i^{\mu_{|S_i|}} | S_i \stackrel{\text{i.i.d.}}{\sim} U(0, \kappa) \text{ for } \mu. \in S_i \quad (26)$$

$$S_i | n_i \sim U([M]^{n_i}) \quad (27)$$

287 The number of active latents n_i are again i.i.d., but there is no constraint on the mean of the distribution
 288 (unlike BatchTopK which constrained the mean of n_i to equal K):

$$n_1, \dots, n_N \stackrel{\text{i.i.d.}}{\sim} P, \quad (28)$$

289 which completes the analysis for JumpReLU SAE. \square

290 **B Stationarity measures**

291 LLM activations are empirically non-stationary across the sequence. We quantify the non-stationary
292 nature by measuring autocorrelations and the U-statistic.

293 **B.1 Autocorrelation**

294 We compute autocorrelation by selecting evenly spaced tokens across the sequence and measuring the
295 cosine similarity between each token and tokens at various lags in the past. Specifically, for tokens at
296 position t , we compute similarities to tokens at $t - w$ where lag w ranges from 5 to 20. This creates a
297 heatmap where rows represent lag offsets and columns represent token positions.

298 For a stationary process, we expect the autocorrelation pattern to remain consistent across time—that
299 is, the relationship between a token and its historical context should be similar regardless of position
300 in the sequence. This would manifest as similar autocorrelation patterns repeating horizontally across
301 token positions. In contrast, for a non-stationary process where representations evolve over time, we
302 expect the autocorrelation patterns to vary systematically across positions, with columns showing
303 different temporal dependency structures as the sequence progresses.

304 **B.2 U-statistic**

305 We measure the effective dimensionality of LLM representations using a U-statistic based on pairwise
306 cosine similarities.

$$\text{U-stat}(t) = \frac{M^2 - M}{\|\mathbf{G}_t\|_F^2 - M} \quad (29)$$

307 where $\|\mathbf{G}_t\|_F^2$ is the squared Frobenius norm of the Gram matrix. This quantity estimates the effective
308 rank $1/\text{tr}(\mathbf{C}_t^2)$, where $\mathbf{C}_t = \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^T]$ is the second moment matrix and \mathbf{x}_t is the activation vector at
309 time t . Under stationarity, U-stat remains constant. When representations evolve over time, U-stat
310 increases systematically as more orthogonal directions become active.

311 **B.3 Surrogate**

312 For both metrics, we compare LLM activations to a known stationary distribution, referred to as
313 surrogate. The surrogate is obtained from observed LLM activation distribution by randomize the
314 phase in Fourier space. It preserves the magnitude of each dimension while randomizing the phases,
315 which creates a stationary process that maintains the same marginal distribution as the original
316 activations.

317 **C Further Results**

318 **C.1 Autocorrelation plots for other LLMs**

319

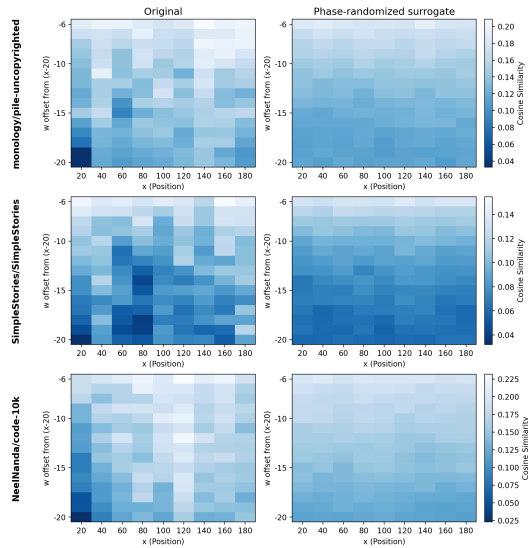


Figure 4: Autocorrelation for Llama 3.1 8B and a surrogate stationary process

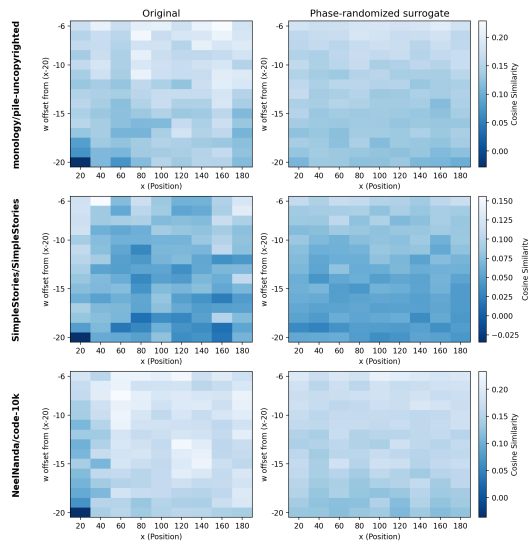


Figure 5: Autocorrelation for Gemma 2 2B and a surrogate stationary process

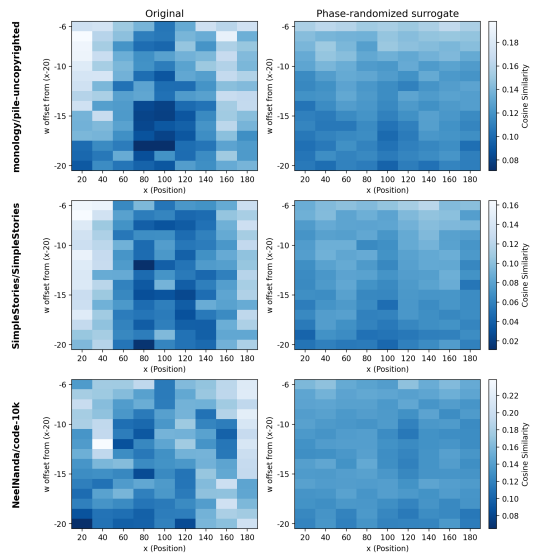


Figure 6: Autocorrelation for GPT 2 small and a surrogate stationary process