

---

# Evaluating Synthetic Activations composed of SAE Latents in GPT-2

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Sparse Auto-Encoders (SAEs) are commonly employed in mechanistic inter-  
2 pretability to decompose the residual stream into monosemantic SAE latents.  
3 Recent work demonstrates that perturbing a model’s activations at an early layer  
4 results in a step-function-like change in the model’s final layer activations. Further-  
5 more, the model’s sensitivity to this perturbation differs between model-generated  
6 (real) activations and random activations. In our study, we assess model sensitivity  
7 to compare real activations to synthetic activations composed of SAE latents. Our  
8 findings indicate that synthetic activations closely resemble real activations when  
9 we control for the sparsity and cosine similarity of the constituent SAE latents.  
10 This suggests that real activations cannot be explained by a simple “bag of SAE  
11 latents” lacking internal structure, and instead suggests that SAE latents possess  
12 significant geometric and statistical properties. Notably, we observe that our syn-  
13 thetic activations exhibit less pronounced activation plateaus compared to those  
14 typically surrounding real activations.

## 15 1 Introduction

16 Neural networks often exhibit polysemanticity, where individual neurons fire for multiple features  
17 Olah et al. [2017]. To explain this, the theory of *superposition* suggests that neural networks represent  
18 more features than they have dimensions, with features linearly represented as directions in activation  
19 space [Elhage et al., 2022, Bricken et al., 2023]. However, the claim that all features are represented  
20 as directions remains speculative [Engels et al., 2024, Smith, 2024, Olah, 2024].

21 Sparse Auto-Encoders (SAEs) have become increasingly popular for decomposing a model’s residual  
22 stream into more interpretable latents Sharkey et al. [2022], Bricken et al. [2023], Cunningham et al.  
23 [2023]. As reliance on SAEs grows, it is crucial to verify that they accurately capture model-used  
24 abstractions.

25 Neural networks employing superposition to represent features must address the challenge of interfer-  
26 ence to maintain performance [Hänni et al., 2024]. This necessitates an ability to accurately extract  
27 individual features while mitigating noise from “nearby” features in the representation space.

28 Heimersheim and Mendel [2024] observed two key phenomena related to this: activation plateaus and  
29 directional sensitivity. These are characterized by changes in the L2 distance of model activations at  
30 the final layer in response to early layer perturbations. Activation plateaus indicate model robustness  
31 to small amounts of noise, while directional sensitivity refers to the model’s varied response to  
32 perturbations in different directions. Importantly, activation plateaus are present around model-  
33 generated activations (real) but not around random points sampled from the Gaussian approximation  
34 of the distribution of model-generated activations (random).

---

\*Equal contribution

35 In this paper, we generate synthetic activations composed of SAE latents and test if they behave  
36 like real activations. We investigate whether arbitrary combinations of SAE latents, “bags of SAE  
37 latents”, can produce activations resembling real ones, and explore the role of latent sparsity and  
38 cosine similarity in this process. Our key contributions include:

- 39 1. We find that the “bag of SAE latents” approach is not sufficient to produce synthetic  
40 activations that resemble model-generated (real) activations.
- 41 2. We find that the sparsity of the top SAE latent, the relative latent activations, and the cosine  
42 similarity between the active latents and the top latent play an important role in determining  
43 whether synthetic activations behave like real activations.
- 44 3. The performance of synthetic activations in the sensitivity experiment does not transfer to  
45 the activation plateau experiment that we conduct. We find that synthetic activations do not  
46 have activation plateaus around them like real activations do.

## 47 **2 Background**

48 Our experiments are based on the setup described in Heimersheim and Mendel [2024], wherein they  
49 perturbed model activations at an early layer and measured the effect it had on the L2 distance of late-  
50 layer activations. They investigated activation plateaus and sensitive directions in GPT-2, motivated  
51 by the error correction mechanism predicted by computation in superposition. They explored  
52 two key predictions: (1) model-generated activations should be resistant to small perturbations,  
53 exhibiting "activation plateaus", and (2) perturbations towards model-generated activations should  
54 affect model output more quickly than towards random directions. Their findings supported both of  
55 their predictions, providing evidence for an error correction mechanism used by the model to suppress  
56 small amounts of noise. This research aimed to better understand computation in superposition and  
57 to find dataset-independent evidence for model features, potentially connecting to SAE research.

## 58 **3 Related Work**

59 Several studies have explored model responses to residual stream perturbations:

60 Janiak et al. [2024] identified stable regions (corresponding to activation plateaus) in the activation  
61 space of transformer-based models, hypothesizing their role in error correction and semantic distinc-  
62 tions. Our work primarily focuses on sensitive directions, though we study activation plateaus around  
63 synthetic activations and compare them against real activations.

64 Gurnee [2024] showed that SAE reconstructions cause larger KL divergence shifts in model outputs  
65 compared to equidistant random vectors when substituted for original activations. While our work  
66 focuses on compositions of SAE latents, we study the effect of SAE reconstruction error on our  
67 experiments (Appendix D).

68 Lee and Heimersheim [2024] investigated SAE reconstruction errors and end-to-end SAE latents,  
69 focusing on individual latent directions. Our work differs by studying compositions of SAE latents.

70 Lindsey [2024] examined the effects of ablating and dampening SAE latents on model performance.  
71 In our study, we focus on composing synthetic activations and studying SAE latent properties.

## 72 **4 Method**

73 We adapt the experimental settings from Heimersheim and Mendel [2024] to test whether synthetic  
74 activations composed of SAE latents exhibit behaviors similar to model-generated (real) activations.  
75 This approach allows us to study key relationships between SAE latents for generating in-distribution  
76 synthetic activations. Section 4.1 outlines the directional sensitivity experiment methodology, Sec-  
77 tions 4.2 and 4.3 describe the activation types tested, and Section 4.4 details the activation plateau  
78 experiment.

79 **4.1 Perturbation Setup**

80 We perturb activations at the last token position in layer 1 (`blocks.1.hook_resid_pre`). The  
 81 unperturbed base activation  $A$  is perturbed towards a direction  $D$ :

$$A_{\text{pert}}(n) = A + 0.5 \cdot n \cdot D$$

82 where  $n$  is the step number (0 to 100), and  $D$  is the normalized difference between base and target  
 83 activations. We use a step size of 0.5, making perturbation norms comparable to typical activation  
 84 norms ( $\simeq 56$ ).

85 We measure L2 distance between original and perturbed activations after the final layer  
 86 (`blocks.11.hook_resid_post`), preferring it over KL divergence for clearer activation plateau  
 87 structure (KL divergence results in Appendix C).

88 To locate blowups, we use the maximum slope (MS) step of the L2 distance curve (Figure 1; we  
 89 discuss alternative metrics in Appendix B).

90 We use GPT2-small [Radford et al., 2019] for our experiments, running inference on random 10-  
 91 token prompts from OpenWebText [Gokaslan and Cohen, 2019]\*. Model-generated activations are  
 92 collected from Layer 1 (`blocks.1.hook_resid_pre`). We employ GPT2-small SAEs [Bloom,  
 93 2024], `sae-lens` [Bloom and Chanin, 2024], and TransformerLens [Nanda and Bloom, 2022] for  
 94 experiments and synthetic activation generation.

95 **4.2 Non-SAE Baselines**

96 In order to compare our setup to previous work [Heimersheim and Mendel, 2024], we run perturba-  
 97 tions towards model-generated (real) and random activations. We sample 1000 prompts and obtain  
 98 base activations, and perturb each base activation in two directions:

- 99 • **Model-generated (real):** Towards a randomly selected activation produced by the model.
- 100 • **Random:** Towards a randomly sampled point from a normal distribution with the same mean  
 101 and covariance as model-generated activations (calculated using 32,000 model-generated  
 102 activations).

103 We plot examples of perturbations towards real and random activations in Figure 1. Both baselines  
 104 have similar base-target distances (mean  $\simeq 40$ , corresponding to step 80). The average cosine  
 105 similarity between model-generated activations (w.r.t. SAE decoder bias) is  $\simeq 0.42$ .

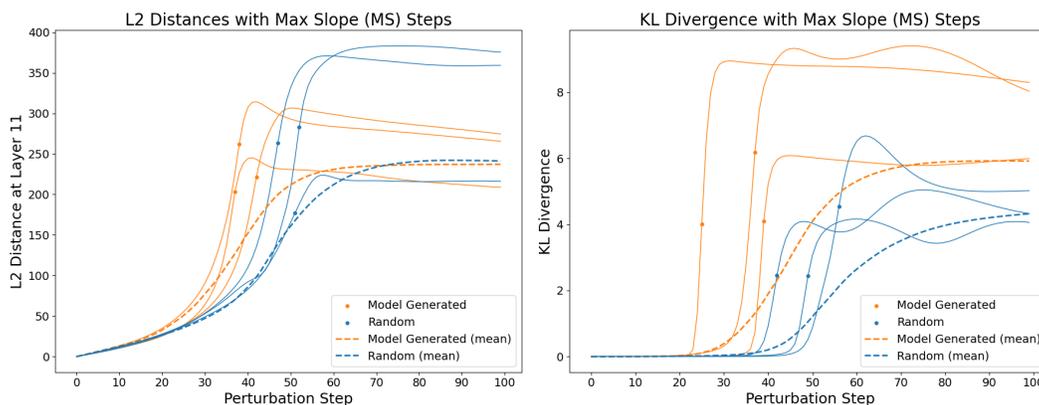


Figure 1: L2 distance (left) and KL divergence (right) between perturbed and unperturbed models for perturbations towards model-generated (orange) and random (blue) activations. X-axis: 100 perturbation steps of 0.5 size each. Dots: maximum slope steps. Dashed lines: average over 1000 perturbations. Initial linear part: activation plateau; sharp rise: blowup.

\*Tokenized dataset link anonymized for review

### 106 4.3 Synthetic Activations

107 We construct synthetic activations using three methods, each incorporating different levels of SAE  
108 latent information:

109 **Synthetic-random:** Randomly selects SAE latents, assigning them base activation’s latent activations.  
110 We present results for this activation type in Appendix A.

111 **Synthetic-baseline:** Accounts for SAE latent sparsities and activations ("bag of SAE latents"):

- 112 1. Encode base activation to obtain active latents and their activations.
- 113 2. Replace each active latent with a random one from the 10 most similarly-sparse latents and  
114 assign the same activation as the original latent.
- 115 3. Decode to obtain synthetic-baseline activation.

116 **Synthetic-structured:** Additionally captures geometric properties of real activations’ SAE latents:

- 117 1. Encode base activation, identifying active latents and their activations. Define `top_base` as  
118 the latent with highest activation.
- 119 2. Create a list of 100 non-dead SAE latents with the most similar sparsity to `top_base`. Out  
120 of the 100 selected latents, select one that has cosine similarity closest to 0.42 (mean cosine  
121 similarity between two real activations w.r.t. the SAE decoder bias) with `top_base`.
- 122 3. This latent becomes the top latent for our synthetic activation (`top_synth`), and we give it  
123 a latent activation value equal to that of `top_base`.
- 124 4. For each remaining active latent in the base activation:
  - 125 (a) Calculate its cosine similarity (`l_top_cos_sim`) with `top_base`.
  - 126 (b) Select a latent (`l_synth`) that has cosine similarity with `top_synth` equal to  
127 `l_top_cos_sim`.
  - 128 (c) Assign `l_synth` a latent activation value equal to that of `l_base`.
- 129 5. Construct a latent activation vector with zeros for all latents except the latents selected above,  
130 and decode it to obtain the **synthetic-structured** activation.

131 We perform 1000 perturbations per synthetic activation type, as described in Section 4.1.

### 132 4.4 Activation Plateaus

133 To test if synthetic activations exhibit activation plateaus like real activations, we use the following  
134 approach:

- 135 1. Initiate perturbations from four base activation types: model-generated, synthetic-baseline,  
136 synthetic-structured, and random (as described in Section 4.2).
- 137 2. Perturb all base types towards random activations (as described in Section 4.2).
- 138 3. Record the activation plateau (AP) step where L2 distance at Layer 11 crosses 20, indicating  
139 plateau flatness.

140 We perform 1000 perturbations per base type, collecting AP size distributions. Larger AP steps  
141 indicate flatter activation plateaus.

## 142 5 Results

143 Synthetic activations behave differently from real and random activations across our two experiments,  
144 suggesting directional sensitivity and activation plateaus point to different properties of SAE latents  
145 in real activations (details in Appendix E). We primarily focus on studying directional sensitivity,  
146 though we also include our findings regarding activation plateaus below.

147 **5.1 Directional Sensitivity**

148 We perturb real activations towards different activation types and study the model’s sensitivity. Figure  
 149 2 and Table 1 show the distributions and statistics of max slope (MS) steps for L2 distance across  
 150 perturbation types.

151 Perturbations towards real activations cause earlier and more localized blowups compared to random  
 152 activations, indicating higher model sensitivity. Synthetic-baseline activations, while not fully  
 153 replicating real activation behavior, outperform random activations. We use the Kolmogorov-Smirnov  
 154 statistic [Smirnov, 1948] to measure distribution similarities.

155 Synthetic-structured activations more closely resemble model-generated activations than synthetic-  
 156 baseline or random activations do. (Figure 2, Table 1). This suggests that relationships between SAE  
 157 latents are important, and that model-generated activations are not approximated well by “bags of  
 158 SAE latents”. Synthetic-random activations perform worse than synthetic-baseline, validating our  
 159 choice of the latter as a stronger baseline (details in Appendix A).

160 To account for varying distances between base and target activations, we also perform perturbations  
 161 with relative step size (Appendix A). This reduces the gap between synthetic-structured and synthetic-  
 162 baseline performance, as synthetic-baseline activations are typically further from base activations and  
 163 thus cause later blowups. It also decreases performance of synthetic-structured.

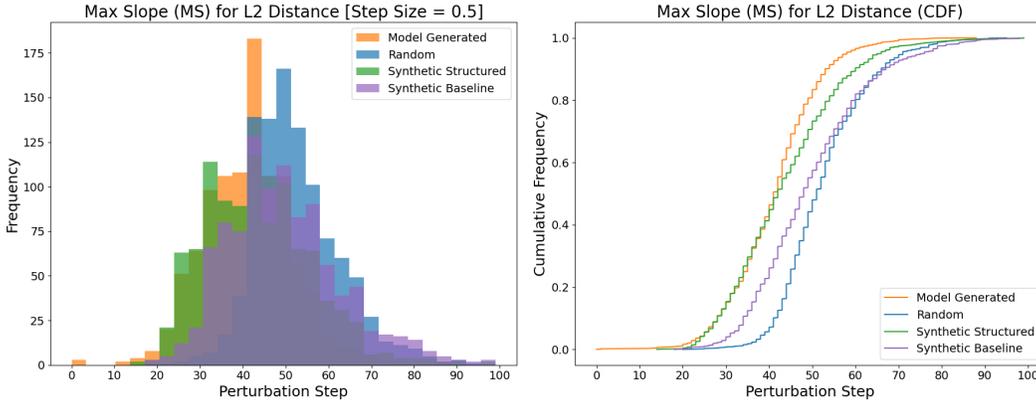


Figure 2: Distribution of MS steps for perturbations towards model-generated (orange), random (blue), synthetic-baseline (purple), and synthetic-structured (green) activations. Left: histogram; Right: cumulative frequency. Synthetic-structured perturbations more closely resemble model-generated ones compared to synthetic-baseline.

Max Slope (MS) step distribution statistics			
Activation Type	Mean	Std dev	KS
Model Generated	41.11	10.40	0.00
Random	52.49	<b>10.21</b>	0.45
Synthetic Baseline	49.61	13.25	0.28
Synthetic Structured	<b>43.48</b>	12.79	<b>0.11</b>

Table 1: Mean, standard deviation, and KS statistic of MS step distributions for perturbations with fixed step size. KS statistic measured against model-generated activations (lower values indicate higher similarity). Synthetic-structured activations most closely resemble model-generated ones.

164 **5.2 Activation Plateaus**

165 Starting from a base of different activation types, we towards random directions to assess their  
 166 activation plateaus. Figure 3 shows the distributions of AP steps for L2 distance across activation  
 167 types.

168 Model-generated activations display pronounced activation plateaus that are not present around in  
 169 random activations. We find that neither synthetic-baseline or synthetic-structured activations show

170 such plateaus, providing further evidence against the "bag of SAE latents" approach but also showing  
 171 that our synthetic-structure activations do not capture all relevant properties of real activations.  
 172 The SAE reconstruction error minimally contributes to the discrepancy between synthetic and model-  
 173 generated activations, we test this in Appendix D.

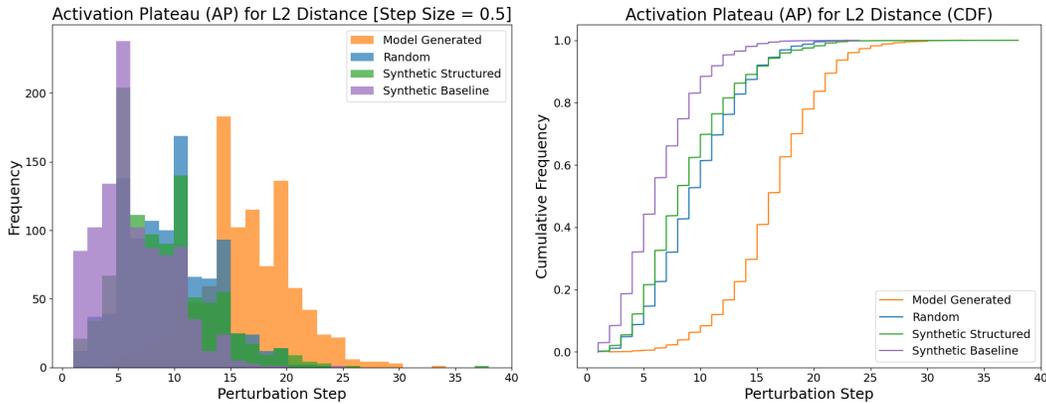


Figure 3: Distribution of activation plateau (AP) steps for perturbations from various activation types towards random activations. Left: histogram; Right: cumulative frequency. Model-generated activations (orange) show flattest plateaus; synthetic-baseline (purple) the steepest. Synthetic-structured (green) and random (blue) activations show similar plateau characteristics.

## 174 6 Limitations

175 The heuristics we use to construct synthetic activations leave room for improvement, as evidenced by  
 176 the gap between them and model-generated activations, especially for activation plateaus. We use  
 177 cosine similarity between SAE latents to capture geometric relationships between them, but leave  
 178 accounting for latent co-occurrence and other relationships between latents for future work.

179 Synthetic activations don't fully match the cosine similarity distribution of SAE latents in model-  
 180 generated activations (Appendix E).

181 Our method leverages information from the base activation in order to construct synthetic activations.  
 182 While this is not ideal, we have verified that using information from a different model-generated  
 183 activation for the construction does not change our results.

184 L2 distance curves' variability may affect our MS metric's effectiveness, as it assumes curve smooth-  
 185 ness (Figure 1). More robust metrics could yield clearer results (Appendix B).

186 Our study is limited to one early layer of GPT2-small and the final token position. Further research  
 187 across different layers, models, SAEs, and context lengths is needed to establish broader applicability  
 188 and generalizability of our findings.

## 189 7 Conclusion

190 Our findings provide additional evidence that GPT-2 is more sensitive to perturbations towards  
 191 model-generated activations than random directions, and that model-generated activations are not  
 192 merely "bags of SAE latents". Leveraging statistical and geometric properties of SAE latents allows  
 193 us to create synthetic-structured activations more similar to model-generated ones, indicating that they  
 194 capture important properties of SAE latents. However, these lack characteristic plateaus of model-  
 195 generated activations, suggesting additional SAE latent properties influence model computation.

196 This presents exciting avenues for future work on model sensitivity to perturbations: developing  
 197 improved synthetic activation construction methods; investigating thresholds for model response  
 198 to latent activation changes; examining model sensitivity to perturbations using interpretable SAE  
 199 latents and contextual information; and analyzing latent ablation-based perturbations to identify key  
 200 contributors to blowups.

201 **References**

- 202 Joseph Bloom and David Chanin. Saelens. <https://github.com/jbloomAus/SAELens>, 2024.
- 203 Joseph Isaac Bloom. Open source sparse autoencoders for all residual stream layers of gpt2-  
204 small, Feb 2024. URL [https://www.alignmentforum.org/posts/f9EgfLSurAiqRjySD/  
205 open-source-sparse-autoencoders-for-all-residual-stream](https://www.alignmentforum.org/posts/f9EgfLSurAiqRjySD/open-source-sparse-autoencoders-for-all-residual-stream).
- 206 Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick  
207 Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec,  
208 Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina  
209 Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and  
210 Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary  
211 learning. *Transformer Circuits Thread*, 2023. [https://transformer-circuits.pub/2023/monosemantic-  
212 features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 213 Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-  
214 coders find highly interpretable features in language models, 2023. URL [https://arxiv.org/  
215 abs/2309.08600](https://arxiv.org/abs/2309.08600).
- 216 Nelson Elhage, Tristan Hume, Catherine Olsson, Neel Nanda, Tom Henighan, Scott Johnston, Sheer  
217 ElShowk, Nicholas Joseph, Nova DasSarma, Ben Mann, Danny Hernandez, Amanda Askell, Kamal  
218 Ndousse, Andy Jones, Dawn Drain, Anna Chen, Yuntao Bai, Deep Ganguli, Liane Lovitt, Zac  
219 Hatfield-Dodds, Jackson Kernion, Tom Conerly, Shauna Kravec, Stanislav Fort, Saurav Kadavath,  
220 Josh Jacobson, Eli Tran-Johnson, Jared Kaplan, Jack Clark, Tom Brown, Sam McCandlish,  
221 Dario Amodei, and Christopher Olah. Softmax linear units. *Transformer Circuits Thread*, 2022.  
222 <https://transformer-circuits.pub/2022/solu/index.html>.
- 223 Joshua Engels, Isaac Liao, Eric J. Michaud, Wes Gurnee, and Max Tegmark. Not all language model  
224 features are linear, 2024. URL <https://arxiv.org/abs/2405.14860>.
- 225 Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. [http://Skylion007.github.io/  
226 OpenWebTextCorpus](http://Skylion007.github.io/OpenWebTextCorpus), 2019.
- 227 Wes Gurnee. Sae reconstruction errors are (empirically) pathological, Mar  
228 2024. URL [https://www.alignmentforum.org/posts/rZPiuFxEsmxCDHe4B/  
229 sae-reconstruction-errors-are-empirically-pathological](https://www.alignmentforum.org/posts/rZPiuFxEsmxCDHe4B/sae-reconstruction-errors-are-empirically-pathological).
- 230 Stefan Heimersheim and Jake Mendel. Activation plateaus and sensitive directions in  
231 gpt2, Jul 2024. URL [https://www.alignmentforum.org/posts/LajDyGyiyX8DNNsuF/  
232 interim-research-report-activation-plateaus-and-sensitive-1](https://www.alignmentforum.org/posts/LajDyGyiyX8DNNsuF/interim-research-report-activation-plateaus-and-sensitive-1).
- 233 Kaarel Hänni, Jake Mendel, Dmitry Vaintrob, and Lawrence Chan. Mathematical models of compu-  
234 tation in superposition, 2024. URL <https://arxiv.org/abs/2408.05451>.
- 235 Jett Janiak, Jacek Karwowski, Chatrik Singh Mangat, Giorgi Giglemiani, Nora Petrova, and Stefan  
236 Heimersheim. Boundaries of stable regions in activation space of llms become sharper with more  
237 compute, September 2024.
- 238 Daniel J. Lee and Stefan Heimersheim. Investigating sensitive directions  
239 in gpt-2: An improved baseline and comparative analysis of saes, Sep  
240 2024. URL [https://www.lesswrong.com/posts/dS5dSgwaDQRoWdTuU/  
241 investigating-sensitive-directions-in-gpt-2-an-improved](https://www.lesswrong.com/posts/dS5dSgwaDQRoWdTuU/investigating-sensitive-directions-in-gpt-2-an-improved).
- 242 Jack Lindsey. How strongly do dictionary learning features influence model behavior? *Transformer  
243 Circuits Thread*, 2024. <https://transformer-circuits.pub/2024/april-update/index.html>.
- 244 Neel Nanda and Joseph Bloom. Transformerlens. [https://github.com/TransformerLensOrg/  
245 TransformerLens](https://github.com/TransformerLensOrg/TransformerLens), 2022.
- 246 Chris Olah. What is a linear representation? what is a multidimensional feature? *Transformer  
247 Circuits Thread*, 2024. <https://transformer-circuits.pub/2024/july-update/index.html>.
- 248 Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization, 2017.

- 249 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
250 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 251 Lee Sharkey, Dan Braun, and Beren Millidge. Taking features out of superposition with sparse autoen-  
252 coders, Dec 2022. URL [https://www.alignmentforum.org/posts/z6QQJbtPkEAX3Aojj/  
253 interim-research-report-taking-features-out-of-superposition](https://www.alignmentforum.org/posts/z6QQJbtPkEAX3Aojj/interim-research-report-taking-features-out-of-superposition).
- 254 N. Smirnov. Table for Estimating the Goodness of Fit of Empirical Distributions. *The Annals  
255 of Mathematical Statistics*, 19(2):279 – 281, 1948. doi: 10.1214/aoms/1177730256. URL  
256 <https://doi.org/10.1214/aoms/1177730256>.
- 257 Lewis Smith. The ‘strong’ feature hypothesis could be wrong, Aug 2024.  
258 URL [https://www.alignmentforum.org/posts/tojtPCCRpKLSHBdpn/  
259 the-strong-feature-hypothesis-could-be-wrong](https://www.alignmentforum.org/posts/tojtPCCRpKLSHBdpn/the-strong-feature-hypothesis-could-be-wrong).

## 260 A Analyzing different perturbations setups and synthetic activations

261 In the main paper we use absolute step size for perturbations, however blowup locations have a  
262 dependence on the distance between the base and target activations, which can make the MS step  
263 distributions with absolute step size misleading. We know that the blowup location does not solely  
264 depend on the distance between the base and target activations, and in order to isolate this property,  
265 we create a distance agnostic setup using relative step sizes. In the relative step size approach, our  
266 perturbations always start at a base activation (A) and end at a target activation (T) using linear  
267 interpolation:

$$A_{\text{pert}}(n) = \left(1 - \frac{n}{100}\right) \cdot A + \frac{n}{100} \cdot T$$

268 where  $n$  is the perturbation step, which goes from 0 to 100. This method ensures that we always  
269 transition from the base activation to the target activation in a fixed number of steps, regardless of  
270 the distance between them. By using relative step size, we remove the dependence of the blowup  
271 location on distance, and instead compare the effect of perturbations purely in terms of the percentage  
272 of base and target activations present at each step. For example, step 50 in this setup implies that the  
273 perturbed activation is made up of 50% base activation and 50% target activation.

274 In the relative step size setup, we find that the MS step distribution for perturbations towards model-  
275 generated activations peaks more strongly around step 50 than in the absolute step size setup. The  
276 blowups are also localized between step 30 and 70, implying that blowups usually happen in the  
277 middle of the perturbation (Figure A.1). We posit that until step 30, the model treats the interpolated  
278 activation as the base activation. This is due to 70% of the interpolated activation coming from the  
279 base activation, and the remaining 30% coming from the target activation being treated as noise. This  
280 effect reverses at step 70, where the model starts treating the interpolated activation as the target  
281 activation, and the 30% that comes from the base activation is considered noise.

282 Our analysis reveals that the MS step distribution for random activation perturbations exhibits  
283 marginally higher variance than the absolute step size setup, with a rightward shift relative to the  
284 distribution for model-generated activation perturbations (Table A.1). This suggests that stronger  
285 perturbations towards random activations are required to induce a blowup compared to model-  
286 generated activations. Furthermore, it indicates that the model is more resilient to random noise than  
287 to noise directed towards another model-generated activation, requiring a greater magnitude of the  
288 former to cause confusion in the model.

289 In this setup, comparing perturbations with synthetic-baseline and synthetic-structured activations  
290 reveals that while synthetic-structured activations still more closely mimic model-generated acti-  
291 vations, the disparity between the two has notably decreased (Figure A.1, Table A.1). This sug-  
292 gests that synthetic-baseline activations less effectively align with the residual stream geometry of  
293 model-generated activations compared to synthetic-structured ones, explaining the latter’s superior  
294 performance in the absolute step size setup. Our findings indicate that considering latent sparsity is  
295 important for synthetic activations to emulate model-generated activations in the relative step size  
296 setup. Consequently, both synthetic-structured and synthetic-baseline outperform synthetic activations  
297 created using the “bag of SAE latents” approach without accounting for sparsity (synthetic-random).

298 We find that when we construct synthetic-structured activations (Section 4.3), omitting the cosine  
 299 similarity constraint on the top latent and instead selecting based on sparsity similarity to the base  
 300 activation’s top latent yields the best-performing synthetic activations in the relative step size setup.  
 301 However, these activations typically have greater distance from the base activation compared to  
 302 synthetic-structured activations. Consequently, their performance in the absolute step size scenario is  
 303 inferior to that of synthetic-structured activations.

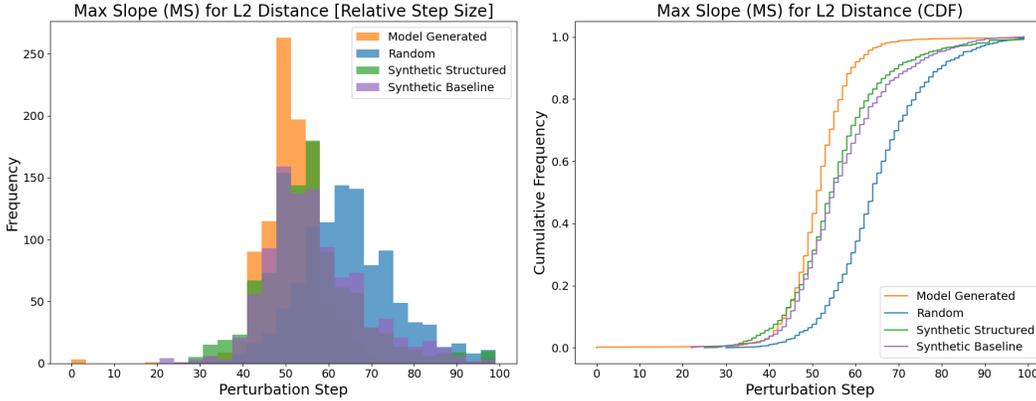


Figure A.1: The distributions of the max slope (MS) steps for perturbations with relative step size towards model-generated (orange), random (blue), synthetic-baseline (purple), and synthetic-structured (green) activations. The left panel shows the counts of MS steps occurring in different bins along the length of the perturbation, and the right panel shows the cumulative frequency for the same. We find that perturbing towards synthetic-structured activations in the relative step size setup is slightly more similar to perturbing towards model-generated activations than perturbing towards synthetic-baseline activations is.

Max Slope (MS) Step Distribution Statistics

Activation Type	Absolute Step Size			Relative Step Size		
	Mean	Std dev	KS	Mean	Std dev	KS
Model Generated	41.21	10.32	0.00	51.65	7.42	0.00
Random	52.49	10.34	0.44	64.89	10.50	0.62
Synthetic Baseline	49.88	12.60	0.31	57.07	11.29	0.27
Synthetic Structured	<b>43.45</b>	12.78	<b>0.11</b>	55.69	11.31	0.22
Synthetic Random	51.30	<b>10.25</b>	0.39	55.25	<b>8.74</b>	0.19
Synthetic Structured (w/o cos sim)	50.17	11.96	0.31	<b>54.47</b>	10.68	<b>0.17</b>

Table A.1: We find that controlling for the sparsity of the top latent and the cosine similarity between the active latents play an important role in making synthetic-structured activations perform well in both absolute and relative step setups. This table contains the mean, standard deviation and KS statistic for MS step distributions for all types of synthetic activations we tested. The KS statistic is measured against perturbations towards model-generated activations, with a lower value meaning higher similarity. The entries in bold show the best match with statistics for model-generated activations.

304 **B Metrics for analysing blowups**

305 In our main analysis, we focus on the maximum slope (MS) as an indicator of the blowup step. In  
 306 this section we share findings using the Area Under the Curve (AUC) and Non Linear (NL) metrics  
 307 to represent important parts of the L2 distance vs perturbation step curve.

308 **B.1 Area Under Curve (AUC)**

309 Our experimental results reveal that certain L2 distance curves deviate from the expected step-  
 310 function-like pattern, causing the MS step to misrepresent the actual blowup location for these curves.  
 311 In contrast, the AUC metric provides a more comprehensive assessment of activation behavior across  
 312 the entire perturbation process. This approach not only identifies the steepest increase point but also  
 313 effectively screens out atypical curves that might otherwise evade detection. AUC calculates the step  
 314 at which the following ratio is maximized:

$$R = \text{area of the triangle defined by } (0,0), (x,0) \text{ and } (f(x),x) / \text{area under the curve } f(x)$$

315 where  $f(x)$  is L2 distance as a function of the perturbation step  $x$ . This method is sensitive to the  
 316 concavity or convexity of the perturbation curve. For predominantly concave curves (where the rate  
 317 of change increases over time), the AUC blowup step tends to occur later, as the triangular area takes  
 318 longer to outpace the actual area under the curve. Conversely, for convex curves (where the rate of  
 319 change decreases over time), the AUC blowup step tends to occur earlier. This property allows the  
 320 AUC method to implicitly capture information about the shape of the perturbation.

321 The AUC metric serves as sanity check, confirming that most perturbations align with expectations.  
 322 Convex L2 distance curves yield early AUC peaks, and Figure B.1 demonstrates that the majority of  
 323 perturbations exhibit the anticipated concave shape. We find that our perturbation results hold for  
 324 AUC step distributions in the absolute step size setup (Table B.1), with structured-synthetic activations  
 325 more closely mimicking model-generated activations compared to synthetic-baseline activations. In  
 326 the relative step size setup (detailed in Appendix A), synthetic-structured and synthetic-baseline  
 327 activations perform similarly. This can be attributed to the higher prevalence of convex curves in  
 328 perturbations towards synthetic-structured activations versus synthetic-baseline activations.

Area Under Curve (AUC) Step Distribution Statistics						
Activation Type	Absolute Step Size			Relative Step Size		
	Mean	Std dev	KS	Mean	Std dev	KS
Model Generated	41.94	11.78	0.00	51.98	9.64	0.00
Random	52.73	13.66	0.43	64.97	16.01	0.59
Synthetic Baseline	49.31	14.20	0.25	56.66	13.20	0.22
Synthetic Structured	43.54	14.99	0.09	54.84	15.51	0.21

Table B.1: We find that our results for the AUC step distributions are similar to those for the MS step distributions. This table contains the mean, standard deviation and KS statistic for AUC step distributions for all the perturbations we perform. The KS statistic is measured against perturbations towards model-generated activations, with a lower value meaning higher similarity.

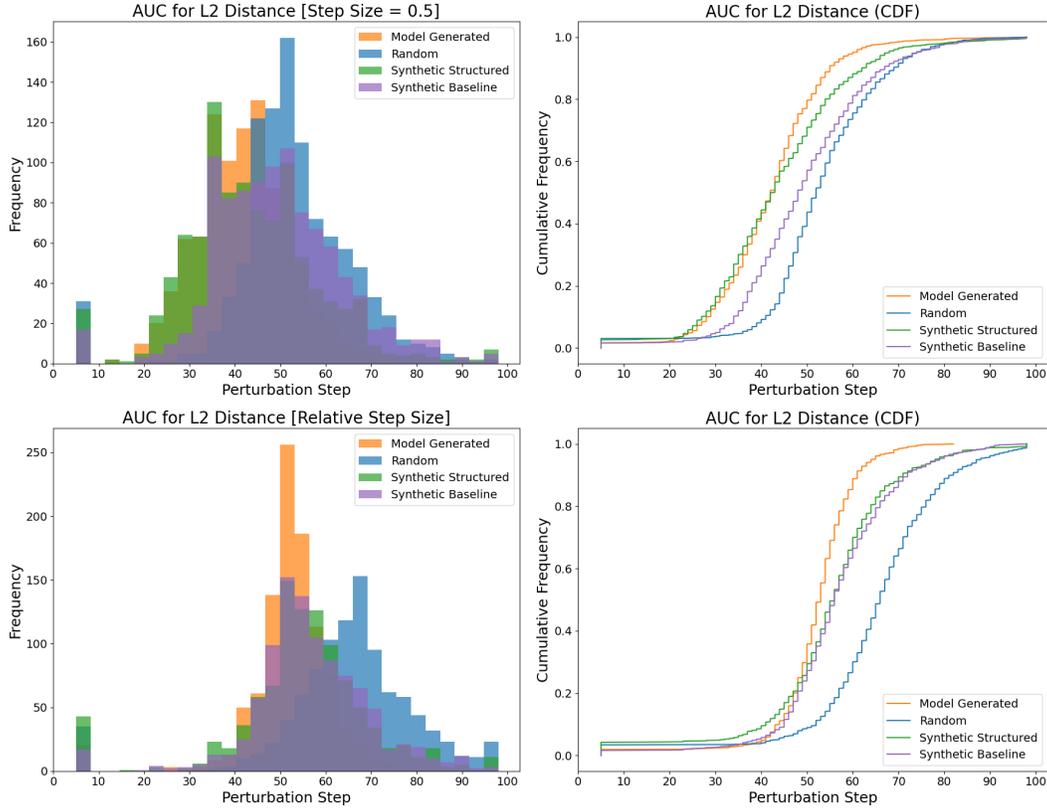


Figure B.1: The distributions of the AUC steps for perturbations with absolute step size (top) and relative step size (bottom) towards model-generated (orange), random (blue), synthetic-baseline (purple), and synthetic-structured (green) activations. The left column shows the counts of AUC steps occurring in different bins along the length of the perturbation, and the right column shows the cumulative frequency for the same. We find that our results for the AUC step distributions are similar to those for the MS step distributions.

## 329 B.2 Non-Linear (NL)

330 Using L2 distance to observe the perturbations reveals that the region before the blowup is not flat,  
 331 but linear with varying slopes (Figure 1). In order to study the size of the initial linear portion of the  
 332 curves, we use the Non-Linear (NL) metric, which points to the earliest step at which the slope of  
 333 the L2 distance vs perturbation step curve deviates from linearity by more than 10% of the initial  
 334 slope. We use this metric as an alternate measure for the size of the activation plateau around the  
 335 base activation along different perturbation directions.

336 We observe that perturbations towards model-generated activations cause the quickest deviation  
 337 from linearity followed by synthetic-structured activations, which is in line with our previous results  
 338 for blowup locations (Figure B.2, Table B.2). However, we find that the deviation from linearity  
 339 occurs the latest during perturbations towards synthetic-baseline activations, which suggests that L2  
 340 distance has a higher initial slope for these perturbations, giving more room for changes in the slope  
 341 before they are classified as a deviation from linearity. In this case, the behavior of synthetic-baseline  
 342 activations provides further evidence that local relationships between SAE latents are important to  
 343 approximate model-generated activations.

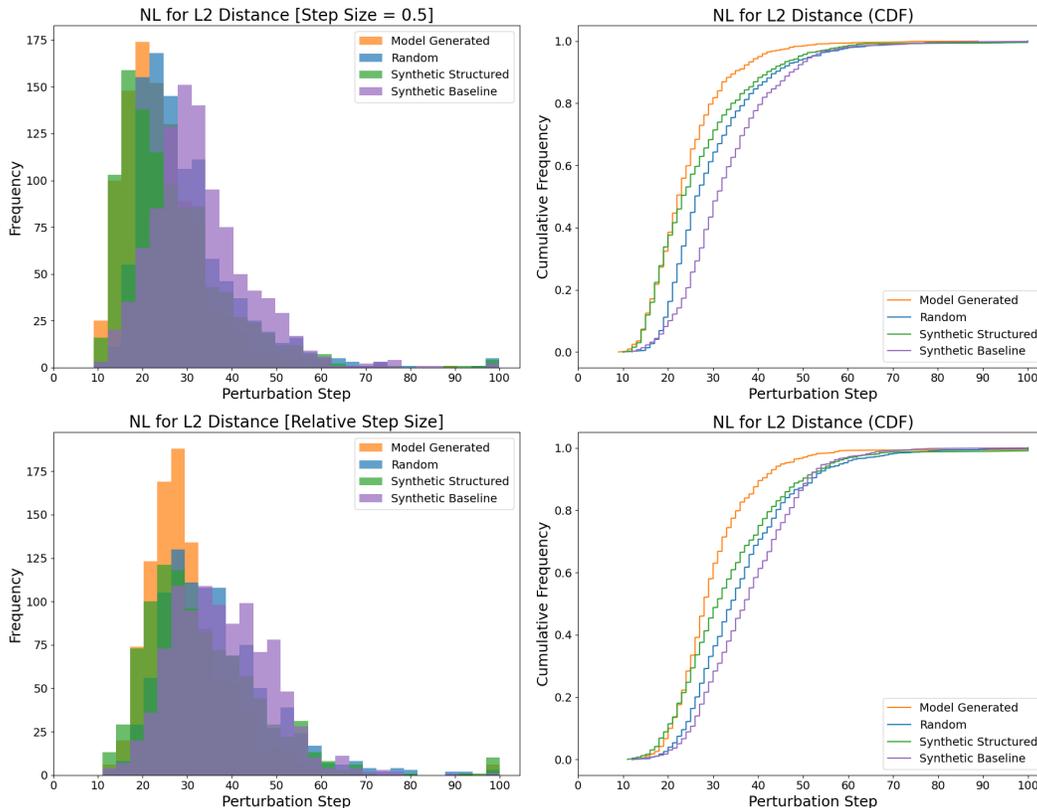


Figure B.2: The distributions of the NL steps for perturbations with absolute step size (top) and relative step size (bottom) towards model-generated (orange), random (blue), synthetic-baseline (purple), and synthetic-structured (green) activations. The left column shows the counts of NL steps occurring in different bins along the length of the perturbation, and the right column shows the cumulative frequency for the same. We find that synthetic-structured activations and random activations behave more like model-generated activations than synthetic-baseline activations do.

Activation Type	Non-Linear (NL) Step Distribution Statistics					
	Absolute Step Size			Relative Step Size		
	Mean	Std dev	KS	Mean	Std dev	KS
Model Generated	24.17	8.87	0.00	29.98	9.95	0.00
Random	29.80	11.72	0.22	36.33	12.33	0.27
Synthetic Baseline	32.90	11.25	0.40	37.91	10.96	0.37
Synthetic Structured	26.72	12.25	0.11	33.69	13.44	0.17

Table B.2: In terms of NL step distributions, we find that synthetic-structured activations perform better than random activations, but synthetic-baseline activations do not. This table contains the mean, standard deviation and KS statistic for NL step distributions for all the perturbations we perform. The KS statistic is measured against perturbations towards model-generated activations, with a lower value meaning higher similarity.

## 344 C KL Divergence

345 While previous works have predominantly used KL divergence as a measure of sensitivity, our  
 346 analysis revealed potential limitations of this approach. We observed that KL divergence produces a  
 347 step-function-like curve even when linear perturbations are performed at the final layer of the model  
 348 right before the unembedding. This behavior suggests that the step-function shape might be an artifact  
 349 of the KL divergence metric itself (or possibly due to softmax), rather than a true representation

350 of activation plateaus. The logarithmic nature of KL divergence may amplify differences as they  
 351 become larger, leading to a more pronounced blowup region and a flatter initial plateau region.

352 With the mentioned caveats in mind, we perform perturbations at Layer 1 and observe their effect on  
 353 KL divergence of the logits distribution instead of L2 distance at Layer 11. Figure C.1 illustrates the  
 354 MS step distribution for KL divergence across different activation types. KL divergence blowups are  
 355 more localized in the relative step size setup than L2 distance blowups, suggesting that the model’s  
 356 output distribution is more robust to noise than the model’s final layer activations, only blowing up  
 357 when more than 40% of the base activation has been replaced. Similar to the results for L2 distance,  
 358 we find that perturbations towards synthetic-structured activations are more similar to perturbations  
 359 towards model-generated activations than synthetic-baseline activations are. The difference between  
 360 synthetic-structured and synthetic-baseline activations is more pronounced for KL divergence than  
 361 L2 distance.

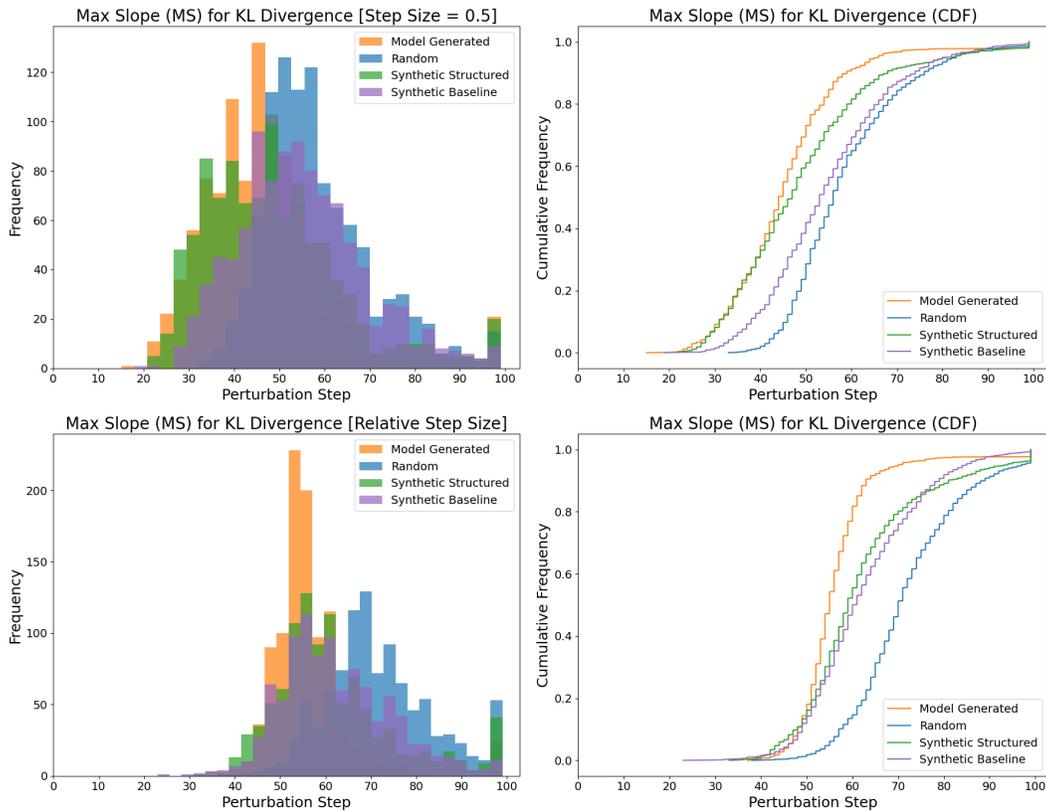


Figure C.1: The distributions of the MS steps for KL divergence of next-token prediction probabilities for perturbations with absolute step size (top) and relative step size (bottom) towards model-generated (orange), random (blue), synthetic-baseline (purple), and synthetic-structured (green) activations. The left column shows the counts of MS steps occurring in different bins along the length of the perturbation, and the right column shows the cumulative frequency for the same. We find that our results for KL divergence are similar to those for L2 distance.

Activation Type	Absolute Step Size			Relative Step Size		
	Mean	Std dev	KS	Mean	Std dev	KS
Model Generated	45.51	12.89	0.00	56.34	9.08	0.00
Random	58.54	<b>12.33</b>	0.47	71.83	<b>11.75</b>	0.69
Synthetic Baseline	54.79	14.02	0.32	62.41	12.05	0.32
Synthetic Structured	<b>48.79</b>	15.64	<b>0.13</b>	<b>61.86</b>	13.51	<b>0.26</b>

Table C.1: We find that our results for KL divergence of next-token prediction probabilities are similar to those for L2 distance at Layer 11. This table contains the mean, standard deviation and KS statistic for MS step distributions for all the perturbations we perform. The KS statistic is measured against perturbations towards model-generated activations, with a lower value meaning higher similarity.

## 362 D Isolating the effect of SAE reconstruction error

363 We denote the reconstruction of an activation  $A$  with  $\text{SAE}(A) = \text{decode}(\text{encode}(A))$ . To isolate  
 364 the effect of SAE reconstruction error on the blowup location, we examine perturbations towards a  
 365 reconstruction of a model-generated target activation  $\text{SAE}(T)$ . We compare these to perturbations  
 366 towards model-generated activations and find that they are very similar, with blowups occurring  
 367 slightly later for perturbations towards SAE reconstructions (Figure D.1, Table D.1). We also find  
 368 that reconstructions of model-generated activations also have plateaus around them. This shows that  
 369 the majority of the difference in our synthetic activations comes from the heuristics we use to select  
 370 latents, and not the SAE reconstruction error.

371 This similarity suggests that SAE reconstructions behave like model-generated activations for the  
 372 most part, and that the reconstruction error causes a small systematic shift in the blowup location.  
 373 This points to some information loss that causes the model to respond slightly less to perturbations  
 374 towards SAE reconstruction, which is relevant for interpreting experiments that use SAE latents.

Activation Type	Absolute Step Size			Relative Step Size		
	Mean	Std dev	KS	Mean	Std dev	KS
Model Generated	41.11	10.40	0.00	51.60	7.82	0.00
Random	52.49	<b>10.21</b>	0.45	65.01	11.19	0.61
SAE Reconstruction	<b>41.49</b>	11.34	<b>0.02</b>	<b>53.34</b>	<b>8.39</b>	<b>0.11</b>

Table D.1: We find that perturbations towards model-generated activations are almost identical to perturbations towards their SAE reconstructions. This table contains the mean, standard deviation and KS statistic for MS step distributions for all the perturbations we perform. The KS statistic is measured against perturbations towards model-generated activations, with a lower value meaning higher similarity.

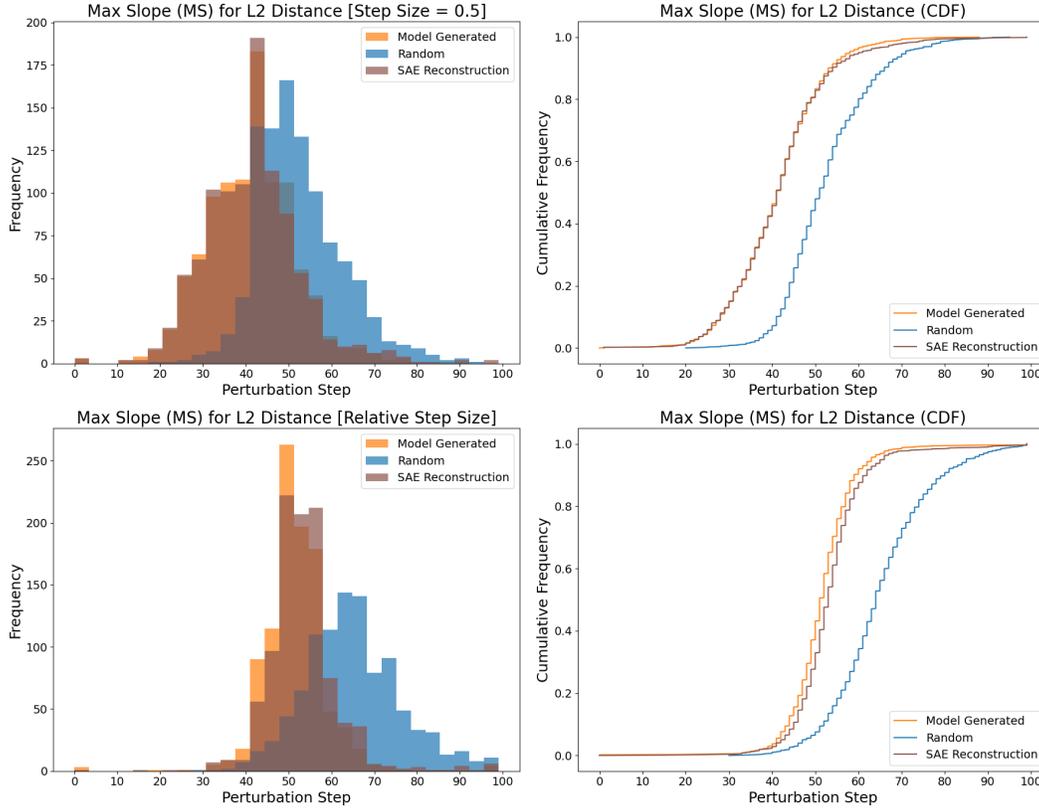


Figure D.1: The distributions of the MS steps for perturbations with absolute step size (top) and relative step size (bottom) towards random activations (blue), model-generated activations (orange), and their SAE reconstructions (brown). The left column shows the counts of MS steps occurring in different bins along the length of the perturbation, and the right column shows the cumulative frequency for the same. We find that perturbations towards model-generated activations and perturbations towards their SAE reconstructions are almost identical.

## 375 E Properties of SAE latents in model activations

376 We observe that model-generated activations with a low SAE reconstruction error contain approx-  
 377 imately 21 active SAE latents on average (Figure E.1 left). The distribution is narrow around the  
 378 mean and falls off very rapidly. The top latent represents around 49% of the total latent  
 379 norm average (Figure E.1 right). The norm falls off rapidly thereafter, with the second top latent  
 380 representing only around 10% on average. The distribution flattens out afterwards where latter ranks  
 381 have similar contribution to the norm.

382 Additionally, we find that model-generated activations are made up of SAE latents that have cosine  
 383 similarity to one another of approximately 0.29 on average (Figure E.2 left), with a distinct peak  
 384 at 0. SAE latents primarily have positive cosine similarity to the top SAE latent, with mean cosine  
 385 similarity of 0.18 (Figure E.2 right) and with a more pronounced peak at 0.

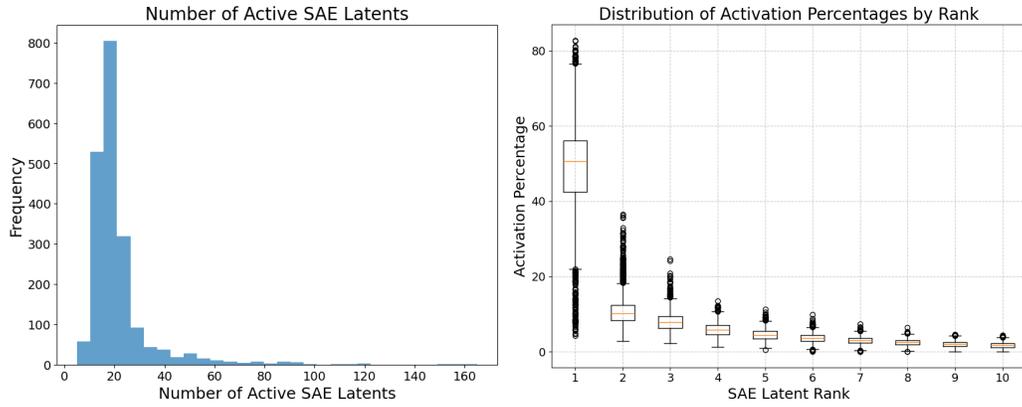


Figure E.1: The distribution of the total number of active SAE latents per activation (left) and the distribution of the percentage of the latent activation norm represented by the top 10 active latents (right) aggregated over 2000 activations.

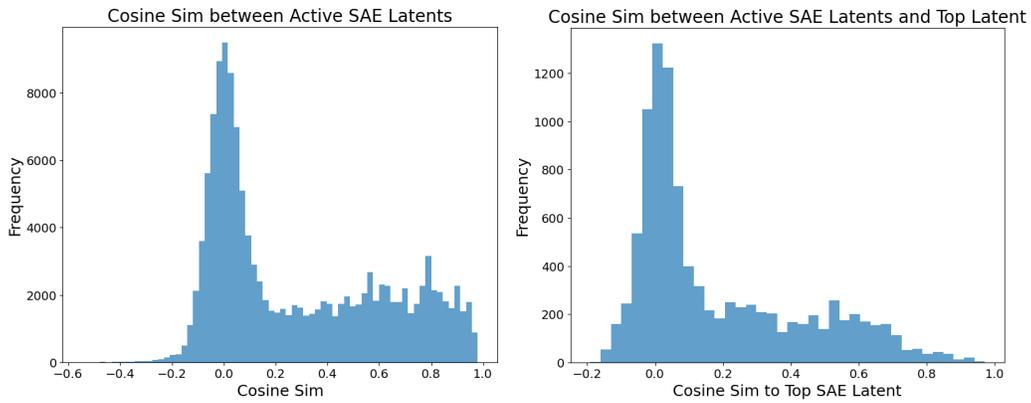


Figure E.2: The distribution of cosine similarities between all active SAE latents per activation (left) and distribution of cosine similarities that active SAE latents have with the top SAE latent (right) aggregated over 2000 activations.