

---

# Decompose Sparsely Where You Should, Absorb Densely Where You Should Not

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Sparse autoencoders (SAEs) are typically trained to reconstruct the **entire** residual  
2 stream through a sparse dictionary, implicitly assuming that all activation content  
3 is amenable to sparse, monosemantic decomposition. We question this assump-  
4 tion and hypothesize that activations contain a low-rank, dense component that is  
5 computationally important to the model yet inherently unsuitable for sparse rep-  
6 resentation, which serves as a major source of the persistent dense latents widely  
7 observed in trained SAEs. To test this, we add a small rank- $r$  linear bottleneck  
8 in parallel with standard SAEs (BatchTopK and Matryoshka), allowing dense  
9 structure to be absorbed before sparse reconstruction. On Gemma-2-2B layer 12,  
10 a rank-24 bottleneck reduces dense latent count by up to 84% while improving  
11 sparse probing and targeted probe perturbation on both architectures at matched  
12 sparsity. The absorbed component is (i) **structurally identifiable** as the top princi-  
13 pal components and outlier dimensions; (ii) **causally necessary**, with removing it  
14 raising next-token cross-entropy by  $7.5\times$ , far exceeding the  $2.8\times$  from removing  
15 the geometrically near-identical top-24 PCA directions; and (iii) **redundantly**  
16 **encoded by sparse dictionaries**, with ablating 787 maximally aligned sparse fea-  
17 tures raising cross-entropy by only  $2.9\times$  and ablating 2,048 topic-aligned features  
18 leaving MMLU topic classification virtually unchanged, whereas removing the  
19 scaffold drops it from 98.7% to chance. Together, our findings identify a compact,  
20 semantically informative and causally important component of residual stream acti-  
21 vations (which we term a **computational scaffold**) that standard sparse dictionaries  
22 represent inefficiently, suggesting that the scope of sparsity-based interpretability  
23 methods warrants careful re-examination.

## 24 1 Introduction

25 Sparse autoencoders (SAEs) have become the predominant tool for extracting interpretable features  
26 from language model activations [3, 8], enabling mechanistic analysis of increasingly large mod-  
27 els [32]. The standard approach trains an SAE to reconstruct the **entire** residual stream activation  
28  $\mathbf{x} \in \mathbb{R}^D$  through a sparse dictionary of learned features. This design carries an implicit assumption  
29 that all meaningful content in the activation can be decomposed into sparse, approximately monose-  
30 mantic directions, an assumption that emerged from how the field has operationalized the Toy Model  
31 of Superposition [11] in practice.

32 Despite widespread adoption of this paradigm, multiple independent lines of evidence suggest the  
33 assumption may be too broad. Sun et al. [31] document persistent high fire rate latents in well-trained  
34 SAEs, classifying them into six functional categories (position tracking, context binding, nullspace,  
35 alphabet, part-of-speech, and PCA reconstruction). Engels et al. [13] identify a linearly predictable  
36 component in SAE residuals that standard architectures systematically fail to capture. Canonical  
37 massive-activation outlier dimensions [1] resist sparse encoding and are known to concentrate in

38 fewer than 20 channels out of thousands. These phenomena have been studied independently as  
39 separate pathologies of SAE training.

40 We hypothesize that these failure modes may reflect a shared structural source, namely activation  
41 content whose computational role is incompatible with sparse dictionary coding at standard scale. If  
42 such content exists, forcing it through a sparse dictionary should waste dictionary capacity on features  
43 that are neither interpretable nor monosemantic, while a dedicated low-dimensional channel should  
44 be able to absorb it, freeing the SAE to specialize on content that is amenable to sparse decomposition.  
45 This hypothesis makes a testable architectural prediction. A small parallel channel, if appropriately  
46 constrained, should simultaneously improve SAE quality metrics **and** reduce the prevalence of dense  
47 latents.

48 To test this prediction, we introduce a **rank- $r$  linear bottleneck** running parallel to standard SAEs.  
49 The SAE operates on the residual  $\mathbf{x} - \text{sg}[\hat{\mathbf{x}}_{\text{dense}}]$ , receiving only what the bottleneck does not capture.  
50 Functional separation is maintained by three design constraints, namely linearity (preventing the  
51 bottleneck from re-encoding sparse features), low rank (preventing SAE starvation), and gradient  
52 isolation (preventing co-adaptation). The bottleneck adds negligible parameters ( $2 \times r \times D$ , less  
53 than 0.7% of an SAE with a 16,384-element dictionary) without requiring changes to the SAE’s  
54 architecture or training procedure.

55 Our contributions are as follows:

- 56 • We propose a simple and effective method for separating activation content that is unsuitable  
57 for sparse decomposition. By placing a low-rank linear bottleneck parallel to a standard  
58 SAE, we allow dense, low-dimensional structure to be absorbed before sparse reconstruction,  
59 without modifying the SAE’s architecture or training procedure. On Gemma-2-2B layer 12,  
60 a rank-24 bottleneck reduces dense latent count by up to 84% and improves sparse probing  
61 and targeted probe perturbation on both architectures at matched sparsity.
- 62 • Through causal intervention and post-hoc analysis, we establish that the absorbed component  
63 is computationally necessary (removing it raises cross-entropy  $7.5\times$ , far exceeding the  $2.8\times$   
64 from removing the geometrically near-identical top-24 PCA directions) and structurally  
65 identifiable as the top principal components and Bondarenko outlier dimensions. Sparse dic-  
66 tionaries encode this content redundantly rather than efficiently, with ablating 787 maximally  
67 aligned sparse features raising cross-entropy by only  $2.9\times$  and ablating 2,048 topic-aligned  
68 features leaving MMLU topic classification virtually unchanged, whereas removing the  
69 component drops it from 98.7% to chance.
- 70 • We provide empirical evidence that the standard practice of training SAEs to reconstruct the  
71 entire residual stream is suboptimal. A subset of activation content is functionally important  
72 for downstream computation yet inherently unsuitable for sparse monosemantic decompo-  
73 sition, which we term a **computational scaffold**. Failing to separate it wastes dictionary  
74 capacity, produces persistent dense latents [31], and degrades downstream extraction quality.  
75 Our results suggest that sparsity-based interpretability methods may have a narrower scope  
76 of applicability than the field’s current practice assumes.

## 77 2 Related work

78 **Sparse autoencoder architectures.** Residual-stream SAEs were established as the standard in-  
79 terpretability tool by Bricken et al. [3], Cunningham et al. [8] and scaled to frontier models by  
80 Templeton et al. [32]. Subsequent work has largely targeted the activation function used to enforce  
81 sparsity, including TopK [16], JumpReLU [29], Gated [28], BatchTopK [4], Matryoshka nested  
82 dictionaries [5], and mixture-of-experts routing [25], together with transcoder variants [10, 27],  
83 end-to-end objectives [2], and large open suites [23, 17] now standardly evaluated by Karvonen  
84 et al. [20]. Our contribution is orthogonal: we leave the SAE’s architecture and training procedure  
85 untouched and add a small parallel linear bottleneck. Whereas Paulo et al. [27] use a full-rank affine  
86 skip on MLP transcoders, we apply a **low-rank** linear projection to **residual-stream SAEs** and show  
87 that the rank constraint is necessary to avoid SAE starvation (Appendix B).

88 **Failure modes of sparse autoencoders.** SAEs exhibit systematic pathologies that suggest structural  
89 limits to the sparse-decomposition assumption. Sun et al. [31] catalogue six classes of dense

90 latents that persist in well-trained SAEs (position, context, nullspace, alphabet, part-of-speech, PCA  
 91 reconstruction), and Engels et al. [13] show that the SAE error vector contains a “dark matter”  
 92 component largely linearly predictable from the input, both indicating that a portion of activation  
 93 content systematically resists sparse decomposition. Complementary failure modes include feature  
 94 absorption and splitting [7], multi-dimensional non-linear features [12], non-canonical and non-  
 95 atomic feature units [22], and seed-dependent decompositions [26]. On the evaluation side, SAEs are  
 96 matched or beaten by simple baselines for probing [19] and steering [33], are not distinguished from  
 97 randomly initialized transformers by standard interpretability metrics [18], and exhibit systematic  
 98 gaps between proxy and downstream metrics [20]. We hypothesise that several of these phenomena,  
 99 including persistent dense latents, dark matter, and the broader resistance of certain content to sparse  
 100 coding, share a common structural source, and provide direct evidence in Section 5.

101 **Structure of language model activations.** Transformer residual streams exhibit pronounced  
 102 structure beyond the feature-superposition picture [11] that motivates SAE design. A long line of  
 103 work documents extreme-magnitude outlier dimensions concentrated in a handful of channels [21,  
 104 9, 1, 30], closely tied to the attention-sink phenomenon [34, 6], while the embedding-geometry  
 105 literature shows that contextual representations are highly anisotropic and dominated by a few top  
 106 principal components [14, 15, 24]. Section 5.1 shows that these outlier and top-PC subspaces overlap  
 107 substantially with the component our bottleneck learns to absorb, connecting the outlier-dimension  
 108 and dense-latent literatures and suggesting both reflect the same low-rank activation structure.

### 109 3 Method

110 We place a linear bottleneck of rank  $r$  parallel to a standard sparse autoencoder (Figure 1). Given  
 111 activation  $\mathbf{x} \in \mathbb{R}^D$ :

$$\mathbf{z} = W_{\text{enc}} \mathbf{x}, \quad W_{\text{enc}} \in \mathbb{R}^{r \times D} \quad (1)$$

$$\hat{\mathbf{x}}_{\text{dense}} = W_{\text{dec}} \mathbf{z}, \quad W_{\text{dec}} \in \mathbb{R}^{D \times r} \quad (2)$$

$$\hat{\mathbf{x}}_{\text{sparse}} = \text{SAE}(\mathbf{x} - \text{sg}[\hat{\mathbf{x}}_{\text{dense}}]) \quad (3)$$

$$\hat{\mathbf{x}} = \hat{\mathbf{x}}_{\text{dense}} + \hat{\mathbf{x}}_{\text{sparse}} \quad (4)$$

112 where  $\text{sg}[\cdot]$  denotes stop-gradient (detach). The SAE receives only the residual after the bottle-  
 113 neck’s contribution, and the bottleneck does not receive gradients through the SAE’s loss. This  
 114 architecture is compatible with any sparse SAE; we validate on BatchTopK ( $|\mathcal{D}|=16384, k=40$ ) and  
 115 MatryoshkaBatchTopK (with the same dictionary size and 5 nested groups).

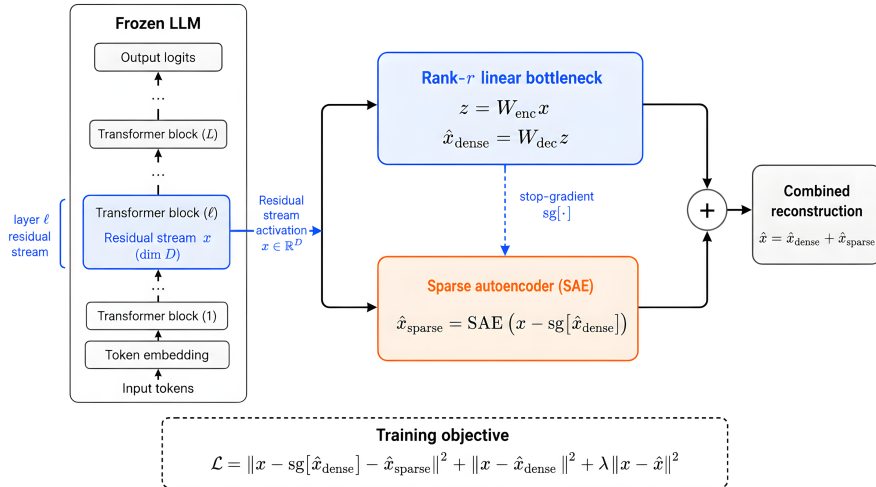


Figure 1: Architecture overview. A rank- $r$  linear bottleneck runs parallel to a standard sparse autoencoder. The SAE operates on the residual  $\mathbf{x} - \text{sg}[\hat{\mathbf{x}}_{\text{dense}}]$  after the bottleneck’s contribution is subtracted with stop-gradient. The combined reconstruction  $\hat{\mathbf{x}} = \hat{\mathbf{x}}_{\text{dense}} + \hat{\mathbf{x}}_{\text{sparse}}$  is trained with a three-term loss (Eq. 5).

Table 1: Effect of adding a rank-24 linear bottleneck to BatchTopK and Matryoshka SAEs. All metrics at strict  $L_0=40$ , matched training conditions.  $\uparrow/\downarrow$  indicate preferred direction. SP, Absorption, and Autointerp are averaged over 3 random seeds. SCR reports direction-1 (spurious-attribute ablation; see Appendix F for metric definition). The bottleneck improves extraction quality metrics (EV, SP, TPP) on both architectures, while intervention based metrics (SCR, RAVEL) show architecture dependent effects analyzed in Appendix F. Autointerp is unchanged on Matryoshka and shows a small ( $-1.2$  pp) regression on BatchTopK that is distributed across fire-rate bins (Appendix G).

Metric	BatchTopK			Matryoshka		
	$r=0$	$r=24$	$\Delta$	$r=0$	$r=24$	$\Delta$
EV $\uparrow$	.793	<b>.805</b>	+0.012	.813	<b>.820</b>	+0.008
SP top-1 $\uparrow$	.749	<b>.758</b>	+0.009	.765	<b>.779</b>	+0.014
SP top-5 $\uparrow$	.851	<b>.864</b>	+0.013	.881	<b>.886</b>	+0.005
TPP@20 $\uparrow$	.023	<b>.054</b>	+0.031	.063	<b>.087</b>	+0.024
Absorption $\downarrow$	<b>.057</b>	.084	+0.027	.011	<b>.009</b>	$-0.002$
SCR@10 $\uparrow$	<b>+213</b>	$-295$	$-508$	+263	<b>+303</b>	+0.040
RAVEL cause $\uparrow$	<b>.582</b>	.558	$-0.024$	.570	<b>.592</b>	+0.022
Autointerp $\uparrow$	<b>.869</b>	.857	$-0.012$	.865	.865	tied
$n_{\text{dense}} \downarrow$	25	<b>4</b>	$-21$	69	<b>44</b>	$-25$
$n_{\text{dead}} \downarrow$	439	<b>305</b>	$-134$	1369	<b>1229</b>	$-140$

116 The training loss combines three terms:

$$\mathcal{L} = \underbrace{\|\mathbf{x}_{\text{sparse\_in}} - \hat{\mathbf{x}}_{\text{sparse}}\|^2}_{\text{sparse reconstruction}} + \underbrace{\|\mathbf{x} - \hat{\mathbf{x}}_{\text{dense}}\|^2}_{\text{bottleneck reconstruction}} + \lambda \underbrace{\|\mathbf{x} - \hat{\mathbf{x}}\|^2}_{\text{full-signal}} \quad (5)$$

117 with  $\lambda=1.0$  and  $\mathbf{x}_{\text{sparse\_in}} = \mathbf{x} - \text{sg}[\hat{\mathbf{x}}_{\text{dense}}]$ . The role of each term, the architectural constraints  
 118 (linearity, low rank, gradient isolation), the rank choice ( $r=24$ ), and control experiments validating  
 119 the linearity and rank constraints are detailed in Appendix B.

## 120 4 Experiments

### 121 4.1 Linear bottleneck improves SAE quality

122 We first ask whether the bottleneck improves SAE quality. We test this by training BatchTopK and  
 123 Matryoshka SAEs with and without a rank-24 bottleneck on Gemma-2-2B (layer 12) at  $L_0=40$ ,  
 124 and evaluating with the SAEbench full suite [5] (training details in Appendix E). Table 1 presents  
 125 the effect of adding a rank-24 linear bottleneck to two SAE architectures under otherwise identical  
 126 training conditions. On both architectures, extraction quality metrics improve consistently. Explained  
 127 variance increases by 0.7–1.2 percentage points, sparse probing top-1 accuracy improves by 0.9–1.4  
 128 points, and targeted probe perturbation (TPP) at threshold 20 roughly doubles. The number of dead  
 129 features also decreases on both architectures, suggesting that absorbing dense content frees dictionary  
 130 capacity.

131 While regressions occur on intervention-based metrics (SCR, RAVEL) and autointerp on BatchTopK,  
 132 these are architecture-dependent and do not appear on Matryoshka. We analyze these effects in  
 133 Appendix F and G.

### 134 4.2 Absorbing the component reduces dense latent count

135 The improvements in Section 4.1 suggest that the bottleneck absorbs content that would otherwise  
 136 burden the sparse dictionary. If our hypothesis is correct, this content should overlap with the  
 137 persistent dense latents documented by Sun et al. [31]. Indeed, the dense latent count (features  
 138 with fire rate  $> 0.1$ ) drops from 25 to 4 at  $r=24$  on BatchTopK and from 69 to 44 on Matryoshka,  
 139 saturating at higher ranks (Figure 2, left).

140 Using the dense latent categories defined by Sun et al. [31], we find that the bottleneck eliminates all  
 141 classifiable dense latents (position tracking, PCA-aligned, and part-of-speech categories) on both  
 142 architectures, with only unclassified dense latents remaining at  $r=24$  (Figure 2, left).

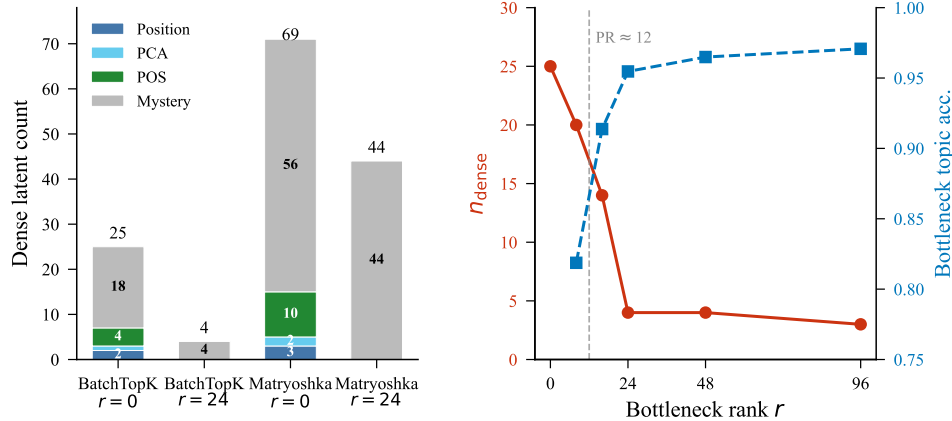


Figure 2: The bottleneck eliminates classifiable dense latents and saturates near the activation’s natural rank. **Left:** dense latent count and category breakdown at  $r=0$  vs.  $r=24$  on both architectures. Classifiable categories (position, PCA, POS) vanish by  $r=24$ ; residual dense latents are unclassified. **Right:** BatchTopK rank sweep ( $r \in \{0, 8, 16, 24, 48, 96\}$ ).  $n_{\text{dense}}$  (left axis, red) drops sharply at  $r=24$  and saturates; MMLU topic accuracy from the bottleneck latent (right axis, blue) increases monotonically and plateaus beyond  $r=24$ . Dashed line marks the participation ratio (PR $\approx 12$ ). Full numerical values in Appendix D.

### 143 4.3 Sweeping bottleneck capacity confirms low-rank structure

144 We further verify whether the absorbed component is low-rank by training a family of BatchTopK  
 145 SAEs across  $r \in \{0, 8, 16, 24, 48, 96, 2304\}$  with all other hyperparameters held fixed. Figure 2  
 146 (right) presents two key metrics across this sweep.

147 Dense latent count ( $n_{\text{dense}}$ , features with fire rate  $> 0.1$ ) drops sharply at  $r=24$  and saturates at  
 148 higher ranks. Topic classification accuracy from the bottleneck latent (a linear probe trained on the  
 149  $r$ -dimensional representation to predict MMLU subject) increases monotonically, reaching 95% at  
 150  $r=24$  with marginal gains beyond. This empirical saturation aligns with the structural rank estimate  
 151 of  $\sim 12$  effective dimensions (Appendix B), so that  $r=24$  provides approximately twice the estimated  
 152 natural rank and both metrics confirm the absorbed component is low-rank.

153 The  $r=2304$  endpoint serves as a starvation control. When the bottleneck has full-rank capacity, it  
 154 absorbs the entire activation (EV=1.00) and the SAE degenerates to noise-level feature selection (SP  
 155 top-1 = .661). This confirms that the low-rank constraint is a necessary design element.

## 156 5 Characterizing the absorbed content as a computational scaffold

157 We now characterize the absorbed content and argue that it constitutes a **computational scaffold**,  
 158 a component that carries coarse contextual information, that the model’s downstream computation  
 159 depends on, and that resists efficient sparse encoding. We establish this through the scaffold’s  
 160 structural identity (Section 5.1), causal necessity (Section 5.2), and redundant representation under  
 161 sparse dictionaries (Section 5.3).

### 162 5.1 The scaffold aligns with top principal components and outlier dimensions

163 To characterize what the bottleneck learns, we measure how much of each principal component’s  
 164 variance is captured by its reconstruction  $\hat{\mathbf{x}}_{\text{dense}}$ . Figure 4 shows a clean step-function pattern, where  
 165 a rank- $r$  bottleneck explains  $>99\%$  of the variance along each of the first  $\sim r$  principal components,  
 166 then drops sharply to the  $r=0$  baseline. The plateau breakpoint tracks the rank, confirming that the  
 167 bottleneck converges to approximately the top- $r$  PCA subspace, as expected from minimizing rank- $r$   
 168 reconstruction MSE.

169 We next ask whether this top-variance subspace corresponds to the outlier dimensions previously  
 170 documented in transformer activations. Per-neuron L2-magnitude rank and top-24 PC loading rank  
 171 correlate at Spearman  $\rho = 0.77$ , with 47 of the top 50 neurons by each measure coinciding (Figure 3,

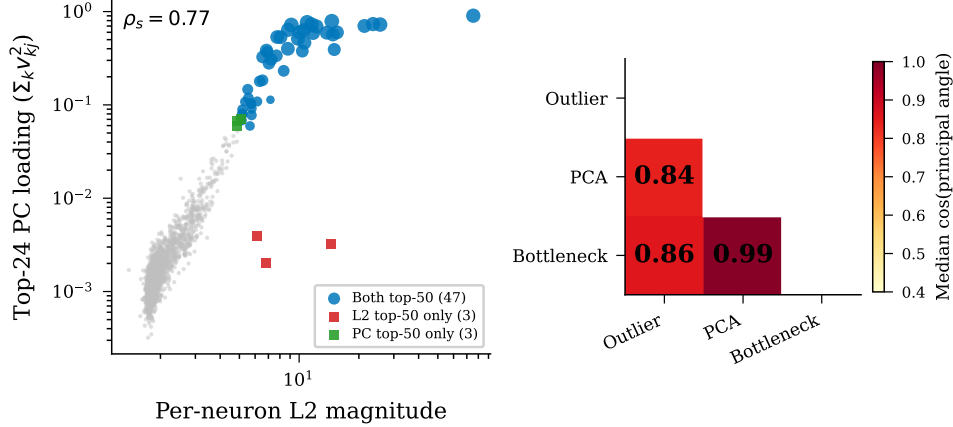


Figure 3: Alignment between outlier dimensions, top principal components, and the learned bottleneck. **Left:** Per-neuron L2 magnitude vs. top-24 PC loading for all 2304 residual-stream neurons (log–log). Blue circles mark the 47 neurons in both the L2-magnitude and PC-loading top-50; dot size is proportional to the bottleneck encoder weight  $\|W_{\text{enc}}[:,j]\|$ . Red and green squares mark the 3 neurons in only one top-50 set. **Right:** Median cosine of principal angles between three 24-dimensional subspace definitions. The bottleneck nearly coincides with the PCA subspace (0.99) and aligns strongly with outlier axes (0.86).

172 left). The dominant-variance directions that the bottleneck absorbs thus largely coincide with the  
 173 canonical outlier dimensions identified by Bondarenko et al. [1].

174 Finally, we compare the three definitions of the 24-dimensional component directly. The subspaces  
 175 defined by outlier axes, top PCA directions, and the learned bottleneck yield pairwise median cosine  
 176 similarities between 0.84 and 0.99 (Figure 3, right). Three operationally distinct procedures, namely  
 177 magnitude ranking, variance decomposition, and end-to-end learning, converge on the same structural  
 178 component of the residual stream.

## 179 5.2 Removing the scaffold disrupts next-token prediction

180 To test the causal importance of this component, we  
 181 project out the bottleneck’s learned 24-dimensional  
 182 subspace from layer-12 activations and measure the  
 183 effect on next-token prediction (full protocol in Ap-  
 184 pendix I). Against a baseline of CE = 1.90, removing  
 185 the component raises CE to 14.3, a  $7.5\times$  increase.  
 186 Two independently trained bottlenecks (BatchTopK  
 187 and Matryoshka, both  $r=24$ ) produce nearly identi-  
 188 cal effects (CE 14.3 and 13.8), confirming that the  
 189 learned component is architecturally stable. Scaling  
 190 the bottleneck contribution by  $\alpha \in [0, 1]$  yields a  
 191 smooth monotonic dose-response (Figure 5, right),  
 192 ruling out a binary on/off role and indicating that  
 193 the model depends continuously on this component’s  
 194 magnitude.

195 Notably, despite the close geometric alignment be-  
 196 tween the bottleneck and the top principal compo-  
 197 nents established in Section 5.1, removing the top-24 PCA directions raises CE to only 5.37 (Figure 5,  
 198 left). This large gap in causal effect suggests that joint training with the SAE steers the bottleneck  
 199 toward computationally critical directions within the top-PC neighborhood that static PCA does not  
 200 isolate.

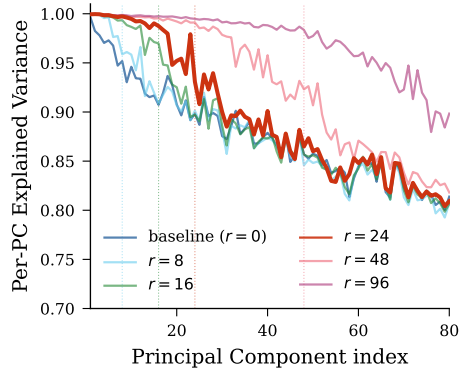


Figure 4: Per-PC explained variance as a function of principal component index.

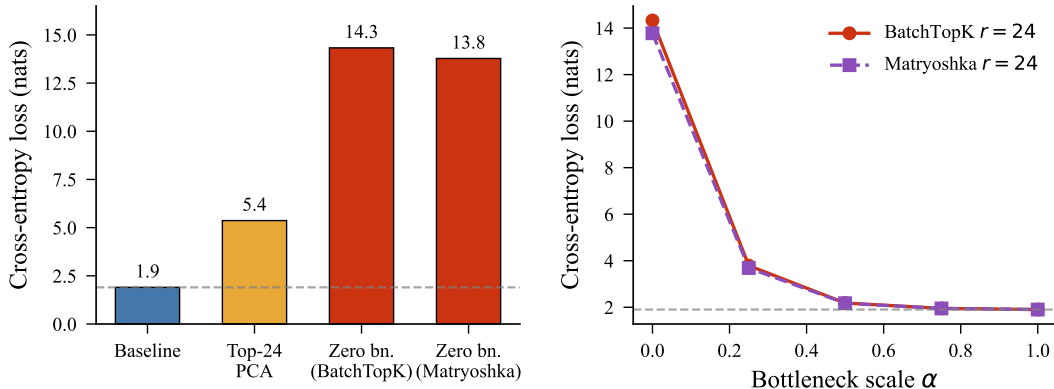


Figure 5: Causal intervention on the bottleneck’s learned 24-d component. **Left:** condition comparison. Removing the bottleneck component (zero) raises CE significantly on both architectures, removing a random 24-d component has negligible effect, and top-24 PCA removal is intermediate. **Right:** dose-response. Scaling the bottleneck contribution  $\alpha \in [0.25, 1.0]$  produces a smooth CE gradient. Dashed lines mark baseline and zero bottleneck reference.

201 **5.3 Sparse encoding of the scaffold is redundant**

202 The results above establish that the scaffold is causally important and geometrically identifiable.  
 203 We now ask whether the sparse dictionary encodes this content efficiently or redundantly. If the  
 204 SAE had concentrated the scaffold into a small number of dedicated features, ablating them should  
 205 approximate the effect of removing the component directly. If instead the scaffold is spread thinly  
 206 across many features that each partially track it, no tractable subset should suffice.

207 We rank all 16,384 SAE features by decoder weight alignment with the bottleneck column space  
 208 and progressively ablate them. Even at  $K=787$ , where the ablated decoder vectors span 96% of the  
 209 scaffold subspace, CE rises to only 5.5 (Figure 6, left), far short of the 14.3 produced by zeroing the  
 210 24-dimensional component itself. Despite near-complete geometric coverage, sparse ablation cannot  
 211 replicate the causal effect, because these features are activation-gated and fire only on matching  
 212 inputs, leaving substantial residual signal in the scaffold directions.

213 The same asymmetry holds for coarse topic information. A linear probe on the bottleneck’s 24-  
 214 dimensional latent achieves 95% accuracy on MMLU subject classification (57 classes, binary OvR),  
 215 approaching the full sparse code’s ceiling of 98%. The multinomial 57-way variant used for the  
 216 ablation curves in Figure 6 (right) reports the same gap on a different scale (98.7% vs. 99.8%), yet  
 217 both protocols agree that 24 bottleneck dimensions recover nearly all of the sparse code’s topic  
 218 information. Ablating bottleneck dimensions by descending  $\eta^2$  produces a sharp cliff at  $K \approx 20$ ,  
 219 reaching chance by  $K=24$ . By contrast, ablating the top 2,048 sparse atoms by  $\eta^2$  leaves accuracy  
 220 above 99%. The same topic-level content is compactly encoded in 24 bottleneck dimensions but  
 221 redundantly distributed across the sparse dictionary.

222 **6 Discussion: The elephant in the interpretability room**

223 **6.1 The scope of sparsity-based interpretability.**

224 A common working assumption behind sparse-dictionary interpretability is that the full residual  
 225 stream should be decomposed into sparse, approximately monosemantic features. Prior work has  
 226 already identified various symptoms suggesting this assumption may be too broad, from persistent  
 227 dense latents [31] to linearly predictable reconstruction error [13], metric-validity concerns [18, 20],  
 228 and feature non-atomicity [22]. These have largely been studied as separate pathologies. Our work  
 229 suggests they may share a common structural source. The component we isolate occupies only  
 230  $\sim 24$  dimensions of the residual stream, yet it accounts for the dominant share of activation variance  
 231 (Section 5.1), proves to be the most causally important structure we can identify (Section 5.2), and  
 232 carries semantically meaningful content (Section 5.3). At the same time, it is the content that sparse  
 233 dictionaries represent least efficiently, with the sparse code tracking it in a distributed and redundant  
 234 fashion that fails to approximate its causal effect (Section 5.3).

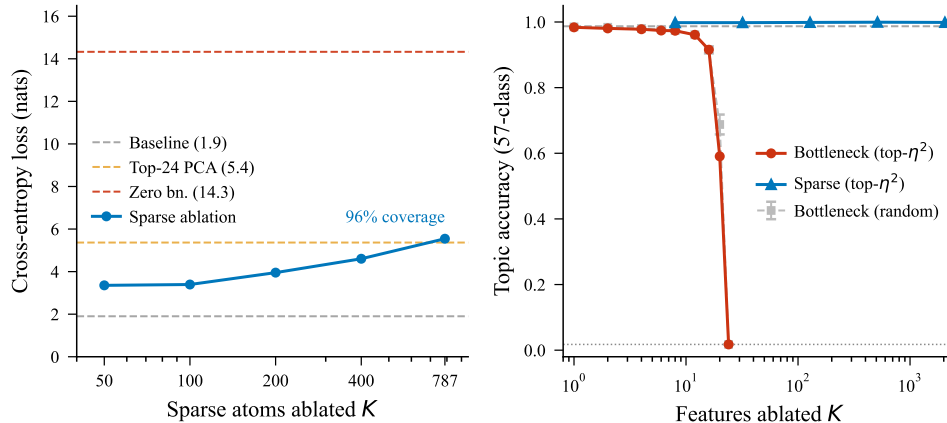


Figure 6: Sparse encoding of the scaffold is redundant. **Left:** CE loss vs. top- $K$  sparse features ablated by decoder  $\alpha$ -alignment with the bottleneck column space. **Right:** topic accuracy (57-class MMLU) under progressive ablation. Bottleneck dim ablation produces a sharp cliff at  $K \approx 20$ , with random ordering tracking top- $\eta^2$  closely (topic is distributed uniformly across the 24 dims).

235 These observations motivate interpreting the bottleneck-absorbed component as a **computational**  
 236 **scaffold**, activation content that is not naturally represented as a small set of sparse monosemantic  
 237 features yet appears functionally important for downstream computation. This interpretation is  
 238 supported by three converging lines of evidence, namely dominant variance share, severe degradation  
 239 upon removal, and semantically structured information content. The top-PC/outlier component we  
 240 isolate is one instance of this category; other forms of scaffold content with different geometric  
 241 profiles may well exist in residual streams (see Open Questions below). The most causally important  
 242 component of the residual stream appears to be the one that the field’s primary interpretability tool  
 243 handles least well. We hope this work serves as an invitation to confront the elephant directly.

## 244 6.2 Implications for evaluation.

245 Current evaluation frameworks for sparse features share an implicit premise, precisely that the  
 246 properties being measured should be isolable to a tractable number of sparse features, which may  
 247 not be the best way of encoding those properties in itself. SCR’s ablation protocol, for instance,  
 248 registers the scaffold’s migration to the bottleneck as a failure of the SAE, while joint ablation of  
 249 bottleneck and sparse components achieves cleaner disentanglement than sparse intervention alone  
 250 (Appendix F).

251 Autointerp exhibits a complementary blind spot. In our cases, autointerp score fails to reflect the  
 252 desirable properties after the bottleneck absorbs persistent dense latents (Appendix G). More broadly,  
 253 Heap et al. [18] demonstrate that autointerp scores fail to distinguish trained transformers from  
 254 randomly initialized ones, suggesting that the metric is insensitive to computational relevance in the  
 255 first place. Both observations suggest that these metrics, while useful, may need to be supplemented  
 256 with evaluations that account for content allocated outside the sparse dictionary.

## 257 6.3 Reinterpreting the success of Matryoshka SAEs.

258 Matryoshka SAEs [5] were designed to address feature absorption through nested prefix dictionaries  
 259 that learn multi-level semantic hierarchies. Our analysis suggests a complementary mechanism  
 260 that may contribute to their success. Their first group  $g_0$  carries a disproportionate share of the  
 261 total explained variance and dense latent population relative to its size, functioning in practice as  
 262 an implicit semi-dense channel for scaffold content (Appendix A). When an explicit bottleneck is  
 263 introduced, it subsumes this role, offloading scaffold content from  $g_0$  and freeing capacity for later  
 264 groups to specialize. The nested prefix design appears to have created, perhaps inadvertently, a  
 265 mechanism for scaffold absorption. The semantic-hierarchy interpretation and the scaffold-absorption  
 266 mechanism are compatible, and may both contribute to Matryoshka’s empirical advantage; the latter,  
 267 however, has not been previously recognized.

## 268 6.4 Open questions.

269 We view our contribution as identifying the first explicit instance of a computational scaffold. This  
270 raises more questions than it answers, and we organize the most pressing ones along three directions.

271 **What else is scaffold?** Our method isolates top-variance content by construction. Content that is  
272 distributed broadly across the principal component spectrum, such as sentiment (Appendix C), falls  
273 outside the bottleneck’s absorption scope. Whether scaffold-like content exists in the middle variance  
274 range, where our current method cannot reach, remains open. The residual dense latents that persist  
275 after bottleneck absorption, unclassified by any test in the taxonomy of Sun et al. [31], hint at scaffold  
276 content beyond the top-variance regime. Similarly, the “dark matter” of Engels et al. [13] is not  
277 substantially reduced by our bottleneck, raising the possibility that it constitutes a distinct form of  
278 scaffold with a different geometric profile.

279 **What computational role does the scaffold play?** The stark causal asymmetry between removing  
280 the scaffold and ablating large numbers of scaffold-aligned sparse features (Section 5.3) suggests that  
281 sparse features may depend on the scaffold as a substrate providing coarse contextual information,  
282 without which fine-grained semantic distinctions cannot be effectively utilized. We emphasize that  
283 this remains an interpretive hypothesis. The scaffold is superficially consistent with the Toy Model of  
284 Superposition [11], which predicts that high-frequency features occupy near-orthogonal directions,  
285 yet its dense, low-rank, collectively encoded character is unlike any feature the TMS framework  
286 considers. Whether it represents an extreme TMS regime or a fundamentally different representational  
287 structure is an important open question.

288 **What is the geometric structure of the scaffold?** The learned bottleneck and the top PCA subspace  
289 are geometrically near-identical yet causally divergent (Section 5.2), indicating that joint training  
290 with the SAE steers the bottleneck toward a causally privileged subset of directions within the  
291 top-PC neighborhood. This gap between geometric overlap and causal effect suggests that the  
292 computationally critical structure within the top-variance subspace is finer-grained than PCA alone  
293 can resolve. Understanding the origin of this discrepancy may shed light on how transformers  
294 organize functionally distinct information within shared variance-dominant subspaces.

## 295 6.5 Limitations.

296 Our experiments use a single model and layer (Gemma-2-2B, layer 12). Because our central claim is  
297 that the working assumption of full sparse decomposability might be too broad, one well-characterized  
298 counterexample suffices to raise the question; nevertheless, cross-model and cross-layer replication  
299 would strengthen confidence in the generality of the architectural prescription. Rank selection ( $r=24$ )  
300 is justified empirically for this setting, and whether the optimal rank scales with model size or layer  
301 depth remains an open question.

302 Our analysis is restricted to residual stream SAEs. Whether MLP transcoders [10, 27] and attention-  
303 output SAEs face analogous scaffold phenomena is unknown. The residual stream aggregates  
304 information from all model components, making it a natural locus for low-rank broadcast structure;  
305 other decomposition targets may exhibit different activation geometry, and the relevance of our  
306 findings to those settings requires independent investigation.

## 307 7 Conclusion

308 We have shown that language model residual streams contain a compact, low-rank component that is  
309 causally critical for downstream computation yet poorly suited to sparse monosemantic decomposition  
310 (a computational scaffold). A simple rank- $r$  linear bottleneck running parallel to a standard SAE  
311 suffices to absorb this component, reducing dense latent counts by up to 84% and improving extraction  
312 quality metrics on two SAE architectures without modifying the SAE itself. The absorbed component  
313 aligns with the top principal components and outlier dimensions, carries coarse semantic content,  
314 and is encoded redundantly rather than efficiently by sparse dictionaries. These findings suggest  
315 that the common practice of training SAEs to reconstruct the entire residual stream conflates two  
316 structurally distinct kinds of activation content, and that explicitly separating them yields better sparse  
317 decompositions. More broadly, our results indicate that the scope of sparsity-based interpretability  
318 warrants re-examination, as the most causally important structure in the residual stream may be  
319 precisely what current methods handle least well.

## References

- 320
- 321 [1] Y. Bondarenko, M. Nagel, and T. Blankevoort. Quantizable transformers: Removing outliers by  
322 helping attention heads do nothing, Nov. 2023.
- 323 [2] D. Braun, J. Taylor, N. Goldowsky-Dill, and L. Sharkey. Identifying functionally important  
324 features with end-to-end sparse dictionary learning, May 2024.
- 325 [3] T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner, C. Anil, C. Deni-  
326 son, A. Aspell, et al. Towards monosemanticity: Decomposing language models with dictionary  
327 learning. *Transformer Circuits Thread*, 2023. URL [https://transformer-circuits.pub/  
328 2023/monosemantic-features](https://transformer-circuits.pub/2023/monosemantic-features).
- 329 [4] B. Bussmann, P. Leask, and N. Nanda. Batchtopk sparse autoencoders, Dec. 2024.
- 330 [5] B. Bussmann, N. Nabeshima, A. Karvonen, and N. Nanda. Learning multi-level features with  
331 matryoshka sparse autoencoders, Mar. 2025.
- 332 [6] N. Cancedda. Spectral filters, dark signals, and attention sinks, Feb. 2024.
- 333 [7] D. Chanin, J. Wilken-Smith, T. Dulka, H. Bhatnagar, S. Golechha, and J. Bloom. A is for  
334 absorption: Studying feature splitting and absorption in sparse autoencoders, Nov. 2025.
- 335 [8] H. Cunningham, A. Ewart, L. Riggs, R. Huben, and L. Sharkey. Sparse autoencoders find highly  
336 interpretable features in language models, Oct. 2023.
- 337 [9] T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer. Llm.int8(): 8-bit matrix multiplication  
338 for transformers at scale, Nov. 2022.
- 339 [10] J. Dunefsky, P. Chlenski, and N. Nanda. Transcoders find interpretable llm feature circuits, Nov.  
340 2024.
- 341 [11] N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds,  
342 R. Lasenby, D. Drain, C. Chen, R. Grosse, S. McCandlish, J. Kaplan, D. Amodei, M. Wattenberg,  
343 and C. Olah. Toy models of superposition, Sept. 2022.
- 344 [12] J. Engels, E. J. Michaud, I. Liao, W. Gurnee, and M. Tegmark. Not all language model features  
345 are one-dimensionally linear, Feb. 2025.
- 346 [13] J. Engels, L. Riggs, and M. Tegmark. Decomposing the dark matter of sparse autoencoders,  
347 Mar. 2025.
- 348 [14] K. Ethayarajh. How contextual are contextualized word representations? comparing the  
349 geometry of bert, elmo, and gpt-2 embeddings, Sept. 2019.
- 350 [15] J. Gao, D. He, X. Tan, T. Qin, L. Wang, and T.-Y. Liu. Representation degeneration problem in  
351 training natural language generation models, July 2019.
- 352 [16] L. Gao, T. D. la Tour, H. Tillman, G. Goh, R. Troll, A. Radford, I. Sutskever, J. Leike, and  
353 J. Wu. Scaling and evaluating sparse autoencoders, June 2024.
- 354 [17] Z. He, W. Shu, X. Ge, L. Chen, J. Wang, Y. Zhou, F. Liu, Q. Guo, X. Huang, Z. Wu, Y.-G.  
355 Jiang, and X. Qiu. Llama scope: Extracting millions of features from llama-3.1-8b with sparse  
356 autoencoders, Oct. 2024.
- 357 [18] T. Heap, T. Lawson, L. Farnik, and L. Aitchison. Automated interpretability metrics do not  
358 distinguish trained and random transformers, Jan. 2026.
- 359 [19] S. Kantamneni, J. Engels, S. Rajamanoharan, M. Tegmark, and N. Nanda. Are sparse autoen-  
360 coders useful? a case study in sparse probing, Feb. 2025.
- 361 [20] A. Karvonen, C. Rager, J. Lin, C. Tigges, J. Bloom, D. Chanin, Y.-T. Lau, E. Farrell, C. Mc-  
362 Dougall, K. Ayonrinde, D. Till, M. Wearden, A. Conmy, S. Marks, and N. Nanda. Saebench:  
363 A comprehensive benchmark for sparse autoencoders in language model interpretability, June  
364 2025.

- 365 [21] O. Kovaleva, S. Kulshreshtha, A. Rogers, and A. Rumshisky. Bert busters: Outlier dimensions  
366 that disrupt transformers, June 2021.
- 367 [22] P. Leask, B. Bussmann, M. Pearce, J. Bloom, C. Tigges, N. A. Moubayed, L. Sharkey, and  
368 N. Nanda. Sparse autoencoders do not find canonical units of analysis, Feb. 2025.
- 369 [23] T. Lieberum, S. Rajamanoharan, A. Conmy, L. Smith, N. Sonnerat, V. Varma, J. Kramár,  
370 A. Dragan, R. Shah, and N. Nanda. Gemma scope: Open sparse autoencoders everywhere all at  
371 once on gemma 2, Aug. 2024.
- 372 [24] J. Mu, S. Bhat, and P. Viswanath. All-but-the-top: Simple and effective postprocessing for word  
373 representations, Mar. 2018.
- 374 [25] A. Mudide, J. Engels, E. J. Michaud, M. Tegmark, and C. S. de Witt. Efficient dictionary  
375 learning with switch sparse autoencoders, June 2025.
- 376 [26] G. Paulo and N. Belrose. Sparse autoencoders trained on the same data learn different features,  
377 Jan. 2025.
- 378 [27] G. Paulo, S. Shabalin, and N. Belrose. Transcoders beat sparse autoencoders for interpretability,  
379 Feb. 2025.
- 380 [28] S. Rajamanoharan, A. Conmy, L. Smith, T. Lieberum, V. Varma, J. Kramár, R. Shah, and  
381 N. Nanda. Improving dictionary learning with gated sparse autoencoders, Apr. 2024.
- 382 [29] S. Rajamanoharan, T. Lieberum, N. Sonnerat, A. Conmy, V. Varma, J. Kramár, and N. Nanda.  
383 Jumping ahead: Improving reconstruction fidelity with jumprelu sparse autoencoders, Aug.  
384 2024.
- 385 [30] M. Sun, X. Chen, J. Z. Kolter, and Z. Liu. Massive activations in large language models, Aug.  
386 2024.
- 387 [31] X. Sun, A. Stolfo, J. Engels, B. Wu, S. Rajamanoharan, M. Sachan, and M. Tegmark. Dense  
388 sae latents are features, not bugs, Nov. 2025.
- 389 [32] A. Templeton, T. Conerly, J. Marcus, J. Lindsey, T. Bricken, B. Chen, A. Pearce, C. Citro,  
390 E. Ameisen, A. Jones, et al. Scaling monosemanticity: Extracting interpretable features from  
391 claude 3 sonnet. *Transformer Circuits Thread*, 2024. URL [https://transformer-circuits.](https://transformer-circuits.pub/2024/scaling-monosemanticity)  
392 [pub/2024/scaling-monosemanticity](https://transformer-circuits.pub/2024/scaling-monosemanticity).
- 393 [33] Z. Wu, A. Arora, A. Geiger, Z. Wang, J. Huang, D. Jurafsky, C. D. Manning, and C. Potts.  
394 Axbench: Steering llms? even simple baselines outperform sparse autoencoders, Mar. 2025.
- 395 [34] G. Xiao, Y. Tian, B. Chen, S. Han, and M. Lewis. Efficient streaming language models with  
396 attention sinks, Apr. 2024.

## 397 **A Matryoshka $g_0$ analysis**

398 Matryoshka SAEs [5] use nested prefix dictionaries to address feature absorption. We observe that  
399 their first group  $g_0$  (512 atoms) functions as an implicit semi-dense channel, carrying 78% of total EV,  
400 harboring 67 of 69 dense latents, and achieving 89.3% topic accuracy at  $k=1$ . Adding the explicit  
401 bottleneck subsumes this role (Figure 7):  $g_0$ 's topic accuracy drops 10 pp (while the bottleneck alone  
402 reaches 95%), its dense latent count drops from 67 to 41, and freed capacity allows later groups to  
403 specialize more effectively.

## 404 **B Design choices and validation**

405 **Three-term loss.** The three terms in Eq. 5 serve complementary roles. The sparse reconstruction  
406 loss trains the SAE to reconstruct the bottleneck's residual  $\mathbf{x}_{\text{sparse\_in}} = \mathbf{x} - \text{sg}[\hat{\mathbf{x}}_{\text{dense}}]$ , and the stop-  
407 gradient ensures this loss does not send gradients to the bottleneck. The bottleneck reconstruction

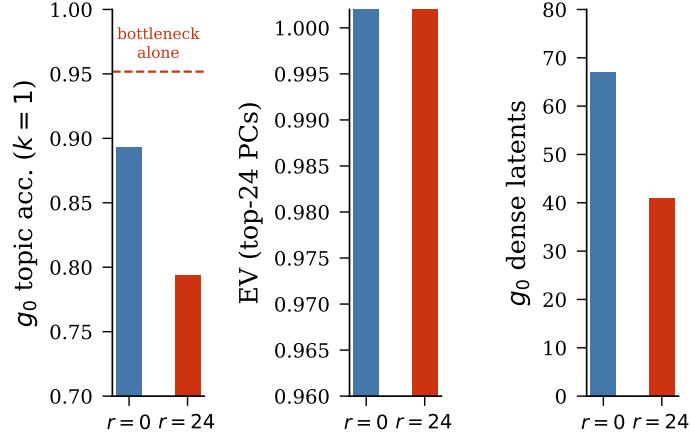


Figure 7: Matryoshka  $g_0$  before and after adding the bottleneck. **Left:**  $g_0$  topic accuracy ( $k=1$ ) drops 10 pp when the bottleneck absorbs topic content (dashed line shows the bottleneck alone achieves 95%). **Center:** per-PC EV on top-24 PCs improves. **Right:**  $g_0$  dense latent count drops from 67 to 41. The explicit bottleneck subsumes the implicit scaffold absorption previously performed by  $g_0$ .

408 loss provides the bottleneck with an independent training signal, driving it toward the minimum-  
 409 MSE rank- $r$  projection. This term is mandatory, as removing it deprives the bottleneck of gradient  
 410 signal and causes degeneration. The full-signal loss ( $\lambda=1.0$ ) provides a joint training signal to both  
 411 components, encouraging coherent combined reconstruction of  $\mathbf{x}$ . This design differs from naive  
 412 joint training ( $\|\mathbf{x} - \hat{\mathbf{x}}\|^2$  alone): under pure joint training, the bottleneck would receive gradients  
 413 only for filling residual gaps left by the SAE, potentially converging to an arbitrary or degenerate  
 414 rank- $r$  complement. The standalone bottleneck loss ensures the bottleneck converges toward the top- $r$   
 415 principal components regardless of the SAE’s behavior, providing a stable functional separation.

416 The four architectural constraints follow, with control experiments validating the linearity and rank  
 417 constraints inline.

418 **Linearity.** A nonlinear bottleneck (GELU activation) can re-encode features the SAE needs,  
 419 defeating the functional separation. Empirically, a GELU-activated bottleneck at rank 24 yields EV =  
 420 .793 (identical to the  $r=0$  baseline), SP top-1 = .721 (below baseline),  $n_{\text{dense}} = 22$  (vs. 4 for linear),  
 421 and bottleneck topic accuracy of 78% (vs. 95% for linear). Classification using the taxonomy of  
 422 Sun et al. [31] confirms the failure, as the GELU variant retains 3 position tracking dense latents  
 423 (vs. 0 for linear at the same rank), indicating that the nonlinear bottleneck does not absorb broadcast  
 424 content. The likely mechanism is that GELU enables the bottleneck to re-encode sparse-feature-like  
 425 content, competing with the SAE rather than complementing it. The fact that SP top-1 **decreases**  
 426 below baseline with GELU (while remaining above baseline with linear) suggests active interference  
 427 between the two components when nonlinearity is present.

428 **Low rank.** A full-rank ( $r=D=2304$ ) linear bottleneck absorbs the entire activation, starving the  
 429 SAE to chance-level performance. The starvation control ( $r=2304$ , full-rank linear) yields EV =  
 430 1.00, leaving the SAE with negligible residual to decompose. Sparse probing drops to .661, RAVEL  
 431 cause to .002, and the SAE’s 40 active features become effectively noise-selected. This confirms that  
 432 the rank constraint is what maintains the cooperative relationship between the two components.

433 **Rank selection.** We set  $r=24$  based on structural observations made **prior to training**. The  
 434 participation ratio of the activation’s per-axis magnitude is 12.3, indicating  $\sim 12$  effective dimensions  
 435 for the heavy-variance structure. Three operationally distinct 24-dimensional subspaces (outlier  
 436 axes, top PCA, and the learned bottleneck) align with pairwise median cosine similarity 0.84–0.99  
 437 (Section 5.1). Per-PC explained variance near-saturates at  $r=24$ .

438 **Gradient isolation.** The stop-gradient in Eq. 3 ensures the SAE optimizes reconstruction of  
 439  $\mathbf{x} - \hat{\mathbf{x}}_{\text{dense}}$  without influencing the bottleneck’s learned component. Without this isolation, the two  
 440 components can co-adapt, undermining the functional separation.

441 Together, these constraints establish a design envelope: the bottleneck must be linear (to prevent  
 442 feature re-encoding) and low-rank (to prevent starvation), with  $r \ll D$  as a hard requirement.

### 443 C Sentiment direction analysis

444 To characterize the boundary of what the bottleneck absorbs, we computed the alignment between a  
 445 CAA-style sentiment direction  $\mathbf{v}_{\text{sentiment}}$  and the top-24 principal components of Gemma-2-2B L12.  
 446 The sentiment direction requires 321 PCs to capture 50% of its norm, and the top-24 PCs hold only  
 447 26.2% of its energy. Accordingly, the bottleneck captures approximately 20% of  $\mathbf{v}_{\text{sentiment}}$ . A linear  
 448 probe on the bottleneck’s 24 latent dimensions for sentiment achieves 60% accuracy (vs. 71% from  
 449 sparse features alone). The bottleneck absorbs **top variance** content. Content that is distributed across  
 450 many principal components falls outside its absorption scope regardless of semantic interpretability.

### 451 D Rank sweep full results

452 Table 2 provides the complete numerical results for the BatchTopK rank sweep summarized in  
 453 Figure 2 (right).

Table 2: Full rank sweep on BatchTopK. Topic = MMLU 57-class accuracy from bottleneck latent only. The  $r=2304$  row is the full-rank starvation control.

Rank	EV	SP top-1	TPP@20	$n_{\text{dense}}$	Topic	EV <sub>top24</sub>	EV <sub>tail500+</sub>	FVU <sub>nl</sub>
0	.793	.749	.023	25	—	.961	.311	.106
8	.797	.759	.081	20	.82	.973	.315	.104
16	.797	.742	.111	14	.91	.985	.317	.102
24	.805	.761	.054	<b>4</b>	.95	.993	.330	.100
48	.805	.771	.063	4	.96	.998	.321	.096
96	.816	.768	.038	3	.97	.999	.327	.090
2304	1.00	.661	.062	45	—	1.00	1.00	.000

### 454 E Training details

455 We train on Gemma-2-2B (base) at layer 12 (`resid_post`) with a dictionary of 16,384 el-  
 456 ements, effective  $L_0=40$  via BatchTopK with strict top- $k$  enforcement at evaluation, using  
 457 `monology/pile-uncopyrighted` for approximately 500M tokens at context length 1024 with  
 458 seed 0. The bottleneck adds  $2rD$  parameters ( $2 \times 24 \times 2304 = 110,592$  for  $r=24$ , or 0.67% of the  
 459 SAE’s parameters).

### 460 F Architecture-dependent SCR effects

461 **Metric definition.** SAEBench’s SCR evaluation computes two directional scores per probe pair.  
 462 Direction 1 ablates features encoding the spurious attribute (e.g., gender) and measures whether  
 463 the target attribute (e.g., profession) prediction recovers toward the unbiased baseline. Direction 2  
 464 performs the reverse ablation. The default `scr_metric` field selects whichever direction has lower  
 465 clean accuracy, which in practice equals direction 2 for most probe pairs in our evaluation. We  
 466 report direction 1 in Table 1 because it directly measures the intended SCR objective (removing  
 467 spurious-attribute features to improve target prediction). Under this definition, the regression on  
 468 BatchTopK is stark ( $+.213 \rightarrow -.295$ ), while the default `scr_metric` would show a mild positive  
 469 change ( $+.239 \rightarrow +.221$ ) that masks the direction-1 collapse.

470 The same rank-24 linear bottleneck produces **opposite** effects on SCR direction-1 depending on the  
 471 host architecture. On BatchTopK, SCR drops from  $+.213$  to  $-.295$ , while on Matryoshka it improves  
 472 from  $+.263$  to  $+.303$ .

473 Disaggregating by probe family reveals the source. On BatchTopK, the regression concentrates in the  
 474 `bias_in_bios` probe (gendered-profession shortcut), where SCR direction-1 improves by +14 to  
 475 +23 percentage points at thresholds 10–50 while SCR metric (averaged over both directions) degrades.

476 The `amazon_reviews` probe (sentiment  $\times$  category) remains positive across ranks. The bottleneck  
 477 absorbs low-rank content that BatchTopK’s sparse features were incidentally capturing alongside  
 478 sparse structure. When this content moves to the bottleneck, sparse feature intervention loses access  
 479 to it, and the SCR metric evaluates this as regression.

480 To understand this dissociation mechanistically, we conduct a dual-attribute ablation on Matryoshka  
 481  $r=24$  using two SCR probe families (Table 3). For `bias_in_bios` (gender  $\times$  profession), the  
 482 bottleneck does not absorb the spurious attribute itself ( $\eta_{\max}^2 = 0.012$ ). Instead, the bottleneck’s  
 483 presence restructures the SAE, concentrating gender information into a small number of atoms (max  
 484  $\eta^2 = 0.120$ , versus 0.020 in the baseline). The small- $K$  ablation protocol then hits these concentrated  
 485 atoms efficiently, producing apparent regression. For `amazon_reviews` (sentiment  $\times$  category),  
 486 the bottleneck partially absorbs sentiment ( $\eta_{\max}^2 = 0.174$ ) and simultaneously spreads the residual  
 487 encoding more diffusely across the SAE (max atom  $\eta^2$  drops from 0.256 to 0.160). The ablation  
 488 finds insufficient concentration to exploit, and the score does not regress.

489 Joint ablation (zeroing the top-8 bottleneck dimensions and top-128 sparse atoms ranked by  $\eta^2$ ) drives  
 490 the gender probe to chance (0.503) while preserving profession accuracy at 0.801, demonstrating that  
 491 the combined bottleneck+SAE system achieves cleaner disentanglement than either component alone  
 492 when both are intervened upon jointly.

Table 3: Dual-attribute ablation on Matryoshka  $r=24$  for two SCR probe families. The bottleneck restructures the SAE differently depending on the spurious attribute’s rank profile. Gender (low-rank) concentrates in sparse atoms, while sentiment (higher-rank) is partially absorbed by the bottleneck and diffused across the sparse code.

Quantity	<code>bias_in_bios</code>	<code>amazon_reviews</code>
Bottleneck $\eta^2$ (spurious) max dim	0.012	0.174
Sparse $\eta^2$ (spurious) max ( $r=0$ )	0.020	0.256
Sparse $\eta^2$ (spurious) max ( $r=24$ )	<b>0.120</b> ( $6\times$ )	0.160 ( $0.6\times$ )
Bottleneck restructures sparse	concentrates	spreads
Joint ablation, target acc.	0.801	0.417
Joint ablation, spurious acc.	<b>0.503</b> (chance)	0.778

493 This pattern has implications for the SCR metric. Its small- $K$  ablation protocol is sensitive to feature  
 494 **concentration** of the spurious attribute, rather than to whether the attribute has been removed from  
 495 the representation. Architectural changes that concentrate rather than redistribute the attribute can  
 496 produce SCR regressions even when the architecture achieves better disentanglement.

## 497 G Autointerp drill-down

498 On BatchTopK, the bottleneck reduces the mean autointerp score from .869 to .857 ( $-1.2$  pp), while  
 499 Matryoshka is unchanged (.865 vs. .865). We run a per-feature drill-down to localize this regression,  
 500 averaging per-latent scores across 3 random seeds (each sampling 1000 features from the same  
 501 dictionary).

502 Table 4 bins features by fire rate and compares mean autointerp score between  $r=0$  and  $r=24$  on  
 503 BatchTopK. The regression is small and distributed across bins rather than concentrated in any single  
 504 fire-rate population. The largest per-bin delta ( $-1.8$  pp) occurs in the  $[0.01, 0.05]$  range. The two  
 505 residual dense latents with fire rate  $> 0.3$  in  $r=24$  score .464, but represent only 2 of 2794 sampled  
 506 features.

507 The score histogram reveals the mechanism. The number of features scoring in the top bin ( $[0.93, 1.0]$ )  
 508 drops from 1428 to 1338 ( $-90$ ), accounting for most of the mean shift. Features newly activated by  
 509 the bottleneck (252 features dead in  $r=0$  but alive in  $r=24$ ) score slightly below the  $r=24$  average  
 510 (.843 vs. .857), contributing a small additional drag. The overall pattern is consistent with the  
 511 bottleneck redistributing dictionary capacity (freeing 134 previously dead slots, shifting fire-rate  
 512 distributions) in a way that marginally reduces the fraction of near-perfect explanation scores without  
 513 degrading any particular feature population.

514 Evaluating only the top-200 features by fire rate yields .699 ( $r=0$ ) vs. .701 ( $r=24$ ), confirming that  
 515 the highest-fire-rate features are not the source of regression. On Matryoshka, the top-200 evaluation

Table 4: Mean autointerp score by fire-rate bin on BatchTopK ( $r=0$  vs.  $r=24$ , 3-seed average). The  $-1.2$  pp overall regression is distributed across bins.  $n$  = number of unique features sampled in that bin.

Fire-rate range	$n_{r=0}$	Score $_{r=0}$	$n_{r=24}$	Score $_{r=24}$	$\Delta$
[0, 0.001)	1587	.915	1384	.911	-.004
[0.001, 0.005)	901	.831	1070	.828	-.004
[0.005, 0.01)	193	.738	253	.737	-.001
[0.01, 0.05)	89	.733	83	.715	-.018
[0.05, 0.1)	6	.583	1	.786	+.202
[0.1, 0.3)	5	.714	1	.786	+.071
[0.3, 1.0]	0	—	2	.464	—

516 drops from .731 to .705, despite the overall random-1000 score being unchanged, suggesting that  
 517 Matryoshka’s group structure redistributes high-fire-rate content differently.

## 518 H Dense latent classification details

519 We reimplemented five of the six automated classifiers described by Sun et al. [31], including position  
 520 tracking (Spearman  $|\rho| > 0.4$  between decoder projection and distance to sentence/paragraph/context  
 521 boundaries), PCA alignment ( $|\cos(\mathbf{W}_{\text{dec}}, \text{PC}_1)| > 0.75$ ), nullspace overlap ( $\alpha_{10} > 0.2$ ), alphabet  
 522 bias ( $\geq 90\%$  shared initial letter among top-100 logit tokens), and meaningful-word part-of-speech  
 523 classification (AUC  $> 0.75$  on Brown Corpus noun/verb/adjective/adverb). We omit context binding,  
 524 which requires LM steering.

525 In the BatchTopK  $r=0$  baseline, these tests classify 7 of 25 dense latents as 2 position tracking, 1  
 526 PCA-aligned, and 4 part-of-speech. Nullspace and alphabet latents are absent, consistent with Sun  
 527 et al. [31]’s finding that these categories emerge predominantly in late layers ( $L \geq 20$ ) rather than  
 528 at L12. In BatchTopK  $r=24$ , only 4 dense latents remain, with 0 classified by the five tests. On  
 529 Matryoshka  $r=0$ , group  $g_0$  (512 atoms) harbors 67 of 69 total dense latents, of which 13 are classified  
 530 (2 position, 2 PCA, 9 part-of-speech). Adding the bottleneck reduces  $g_0$  to 41 dense latents with 0  
 531 classified.

532 To verify that these dense latents track the same structure the bottleneck absorbs, we examine where  
 533 their decoder weights project in PC space. In the BatchTopK  $r=0$  baseline, the 25 dense latents  
 534 project predominantly onto the top-50 principal components, with a median top-PC index of 8 and 13  
 535 of 25 concentrated in the top-10 PCs. In BatchTopK  $r=24$ , only 4 dense latents remain, and 3 of  
 536 these project onto PC index 0 with max cosine  $> 0.70$ , suggesting that these residual dense features  
 537 track the very top of the variance spectrum that the rank-24 bottleneck absorbs almost but not entirely.

538 The connection between dense latents and the scaffold component is tight. Of the 25 dense latents  
 539 in BatchTopK  $r=0$ , 25 of 25 fall in the top decile of decoder weight alignment with the scaffold  
 540 component ( $\alpha > 0.3$ ), and 19 of 25 have  $\alpha > 0.7$ . These features fire at rates of 10–89%, consistent  
 541 with tracking a component present. When counting features with high decoder weight alignment to  
 542 the top-24 outlier-dimension component ( $\alpha_{\text{outlier}} > 0.3$ ), the count drops from 787 in BatchTopK  
 543  $r=0$  to 99 in BatchTopK  $r=24$ , indicating that the bottleneck relieves approximately 688 dictionary  
 544 slots from tracking this low-rank content.

## 545 I Causal intervention protocol

546 All causal interventions in Section 5.2 use 200 sequences (1024 tokens each) drawn from  
 547 monology/pile-uncopyrighted. For each sequence, we run a forward pass through Gemma-  
 548 2-2B and intervene on the layer-12 residual stream activation before continuing to the model head.  
 549 Cross-entropy is computed over all next-token predictions in each sequence and averaged across the  
 550 evaluation set. The baseline (unmodified activations) achieves CE = 1.90.