
ViSAEBench: Cross-Backbone Evaluation of Vision Sparse Autoencoders Reveals Backbone-Dominated Variance and Metric Dissociations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Sparse autoencoders (SAE) are increasingly used to interpret Vision Transformer
2 features, but unlike the language setting, there is no standardized protocol for com-
3 paring vision SAEs and no systematic characterization of how SAE quality depends
4 on the pretrained backbone. We introduce ViSAEBench, a unified evaluation suite
5 covering seven metrics across four interpretability dimensions, including a novel
6 spatial coherence metric specific to vision. Using ViSAEBench, we conduct the
7 first controlled cross-backbone study of vision SAEs: 60 SAEs trained on identical
8 ImageNet-1K activations from five ViT-B backbones spanning four pretraining
9 paradigms. Our central finding is that the choice of pretrained backbone dominates
10 vision SAE behavior more than SAE hyperparameters. A variance decomposition
11 shows that backbone explains over 90% of variance on three metrics and over 60%
12 on five of seven, while SAE hyperparameters dominate only reconstruction error.
13 The starkest instance is categorical: across all configurations tested, SAEs trained
14 on Masked Autoencoder features show no spatial structure beyond chance, while
15 the other four backbones produce strongly spatially structured features. Single-
16 backbone vision SAE evaluations are therefore often measuring properties of the
17 backbone more than properties of the SAE. We further identify two metric-level dis-
18 sociations with practical consequences. First, reconstruction error and downstream
19 task preservation substantially diverge across backbones (Spearman $\rho = -0.70$),
20 so reconstruction error alone cannot be used to compare vision SAEs. Second,
21 monosemanticity, a central SAE quality criterion in language work, does not pre-
22 dict fine-grained classification, indicating that within-feature consistency does not
23 capture the between-class separability downstream tasks require. We release all 60
24 SAE checkpoints and the ViSAEBench evaluation library.

25 1 Introduction

26 Sparse autoencoders (SAEs) have become a standard tool for decomposing neural network repre-
27 sentations into interpretable features [Bricken et al., 2023, Templeton et al., 2024]. For language
28 models, this methodology is supported by unified benchmarks such as SAEBench [Karvonen et al.,
29 2025] and large-scale releases such as Gemma Scope [Lieberum et al., 2024]. Vision SAEs are
30 less standardized: recent work shows that SAEs trained on Vision Transformer (ViT) features can
31 recover monosemantic concepts [Pach et al., 2025] and discrete visual primitives [Sharkey et al.,
32 2025], but studies differ in backbone, metrics, and evaluation datasets. Existing vision representation
33 benchmarks such as VTAB [Zhai et al., 2020] and CLIP probing protocols [Radford et al., 2021]
34 evaluate downstream transfer of the raw backbone, not properties of the SAE-decomposed feature
35 space. Moreover, language SAE metrics such as Fraction of Variance Unexplained (FVU), sparsity,

36 and monosemanticity do not capture the spatial structure unique to ViT patch features: a feature firing
37 on a contiguous image region behaves differently from one diffuse across the patch grid.

38 We introduce **ViSAEBench**, an evaluation suite for SAEs trained on ViT patch features. ViSAEBench
39 spans four dimensions- spatial coherence, reconstruction, concept detection, and disentanglement
40 through seven metrics, including a vision-specific spatial coherence metric (M1) that quantifies
41 whether SAE features activate in spatially contiguous image regions. We apply ViSAEBench to 60
42 BatchTopK SAEs trained on identical ImageNet-1K activations from five ViT-B backbones spanning
43 four pretraining paradigms: DINOv2 (self-distillation) [Oquab et al., 2024], CLIP [Radford et al.,
44 2021] and SigLIP [Zhai et al., 2023] (vision-language contrastive), MAE (masked reconstruction) [He
45 et al., 2022], and DeiT (supervised distillation) [Touvron et al., 2021].

46 Our contributions are: (i) ViSAEBench, a seven-metric evaluation suite for vision SAEs spanning
47 spatial coherence, reconstruction, concept detection, and disentanglement; (ii) an evaluation of 60
48 BatchTopK SAEs trained on identical ImageNet-1K activations from five ViT-B backbones, showing
49 that backbone explains over 90% of variance on three metrics and over 60% on five of seven; (iii) a
50 categorical spatial-coherence separation, with MAE-derived SAEs remaining at the spatial-random
51 null while DINOv2, CLIP, SigLIP, and DeiT produce spatially structured SAE features; and (iv) two
52 metric dissociations, where FVU does not reliably predict downstream preservation across backbones
53 and monosemanticity does not predict CUB-200 fine-grained classification.

54 **2 Related Work**

55 **2.1 SAEs for language model interpretability**

56 Sparse autoencoders have emerged as a central tool for decomposing language model represen-
57 tations into interpretable features [Bricken et al., 2023, Templeton et al., 2024, Lieberum et al.,
58 2024]. Large-scale releases such as Gemma Scope have made SAE features a shared substrate for
59 mechanistic interpretability research. Recent work has extended SAEs to cross-modal settings [Gu
60 et al., 2026, Zhao et al., 2026], exploring vision-language alignment through SAE decomposition.
61 SAEBench [Karvonen et al., 2025] introduced a unified evaluation suite for language SAEs, providing
62 the methodological template we adapt here. Monosemanticity, operationalized as within-feature
63 activation consistency, is treated as a central quality criterion across this line of work; whether it tracks
64 downstream task utility has not been systematically tested. ViSAEBench is the vision counterpart,
65 not a replacement: the metric families overlap (sparsity, reconstruction, probing) but the dimensions
66 specific to vision, particularly spatial coherence, have no language analog.

67 **2.2 SAEs for vision**

68 Recent work has applied SAEs to ViT features and shown that monosemantic visual concepts can
69 be recovered from individual SAE features [Pach et al., 2025]. Earlier patch-feature decomposition
70 work and the ViT-Prisma library [Joseph et al., 2025] established the infrastructure for training
71 SAEs on ViT activations. Parallel efforts have explored alternative architectural surrogates for vision
72 interpretability, including cross-layer transcoders [Chatzoudis et al., 2026], which offer dense global
73 decompositions but lack the sparse, localized structure that enables fine-grained concept isolation.
74 Existing vision SAE work has concentrated on a single backbone per paper and a small number of
75 metrics, leaving cross-architectural comparison undone. ViSAEBench’s contribution relative to this
76 line of work is a controlled cross-backbone study spanning four pretraining paradigms, a unified
77 protocol covering seven metrics across four interpretability dimensions, and a spatial-coherence
78 metric (M1) absent from prior vision SAE evaluation.

79 **2.3 Vision representation benchmarks and spatial structure**

80 Vision representation benchmarks such as VTAB [Zhai et al., 2020], CLIP probing protocols [Radford
81 et al., 2021], and DINO-style linear-evaluation suites [Oquab et al., 2024] evaluate properties of
82 the raw representation through frozen-backbone probing. ViSAEBench evaluates properties of the
83 SAE-decomposed representation, a different question: a backbone with strong probing accuracy
84 may still produce SAE features that are entangled, spatially diffuse, or poorly disentangled. Spatial
85 structure in vision representations has been studied through object-centric learning [Rubinstein
86 et al., 2025], ViT attention probing [Caron et al., 2021, Darcet et al., 2024], and attribution-based

87 localization [Han et al., 2025]. M1 imports a different tradition: Moran’s I [Moran, 1950] is the
 88 standard spatial-autocorrelation statistic in spatial statistics, but to our knowledge has not been
 89 used for evaluating SAE features. We adopt it because it provides a principled null model and a
 90 single-number per-feature score that aggregates cleanly across SAE configurations.

91 3 Background and Setup

92 **BatchTopK Sparse Autoencoders.** A sparse autoencoder is an encoder-decoder pair (E, D) trained
 93 to reconstruct an input activation $\mathbf{a} \in \mathbb{R}^D$ through a sparse intermediate code $\mathbf{c} = E(\mathbf{a}) \in \mathbb{R}^F$,
 94 with $F \gg D$. The TopK variant [Makhzani and Frey, 2014, Gao et al., 2024] enforces sparsity by
 95 retaining only the K largest activations in \mathbf{c} per sample. BatchTopK [Bussmann et al., 2024] relaxes
 96 the per-sample constraint to a per-batch one: in a batch of B activations, the top $B \cdot K$ activations
 97 across the entire batch are retained, allowing simpler samples to use fewer features and complex
 98 samples to use more. We adopt BatchTopK because visual complexity varies substantially across
 99 natural images, and per-sample top- K either underfits complex scenes or wastes capacity on simple
 100 ones. The training objective is

$$\mathcal{L} = \mathbb{E}_{\mathbf{a}} [\|\mathbf{a} - D(\text{BatchTopK}_K(E(\mathbf{a})))\|_2^2]. \quad (1)$$

101 To verify our cross-backbone findings are not specific to BatchTopK, we replicate the full evaluation
 102 on JumpReLU SAEs across all five backbones; cross-family rank correlations are at least 0.80 on
 103 every metric, with six of seven significant at $p < 0.05$ (Appendix Q).

104 **ViT backbones evaluated.** We evaluated five ViT-Base backbones that span four pretraining
 105 paradigms: DINOv2 (self-distillation [Oquab et al., 2024]), MAE (masked patch reconstruction [He
 106 et al., 2022]), CLIP and SigLIP (vision-language contrastive [Radford et al., 2021, Zhai et al., 2023]),
 107 and DeiT (supervised classification with distillation [Touvron et al., 2021]). All five share the ViT-B
 108 architecture (12 transformer blocks, 768 hidden dimension) and use 16×16 patches, except DINOv2,
 109 which uses 14×14 patches. We tap activations from block 11 of each backbone, near the top of the
 110 residual stream where features have been observed to be most semantically discriminative. Sensitivity
 111 to layer choice is reported in Appendix D.

112 **SAE training protocol.** We train 12 SAE configurations per backbone, sweeping expansion
 113 factors $\{8, 16, 32\}$ and sparsity levels $K \in \{64, 128, 192, 256\}$ on 1.28M ImageNet-1K training-split
 114 activations cached from each backbone’s block 11. All SAEs are trained with identical optimizer,
 115 learning rate, and step counts across backbones; full hyperparameters are reported in Appendix A.
 116 The full sweep yields 60 SAEs (5 backbones \times 3 expansion factors \times 4 K values), all released on
 117 HuggingFace under the `visaebench` organization.

118 4 The ViSAEBench Evaluation Suite

119 4.1 Spatial Coherence: Feature Localization Score (M1)

120 **Motivation.** A vision SAE feature is interpretable in the localization sense if it activates on a
 121 spatially contiguous region of the input image. Existing SAE metrics do not capture this: FVU
 122 measures activation-space reconstruction without reference to spatial layout, sparse probing measures
 123 concept presence without regard to where the concept activates, and monosemanticity is invariant to
 124 spatial structure within any single image.

125 **Definition.** Let $\mathbf{f}_i \in \mathbb{R}^P$ denote the activation vector of SAE feature i across the P patches of one
 126 image, where $P=196$ for patch-16 backbones at 224^2 input and $P=256$ for DINOv2’s patch-14
 127 backbone at 224^2 input. We define the per-feature, per-image Moran’s I [Moran, 1950] as

$$I(\mathbf{f}_i) = \frac{P}{\sum_{p,q} w_{pq}} \cdot \frac{\sum_{p,q} w_{pq} (f_{i,p} - \bar{f}_i)(f_{i,q} - \bar{f}_i)}{\sum_p (f_{i,p} - \bar{f}_i)^2}, \quad (2)$$

128 with $w_{pq}=1$ for 8-connected (queen-contiguity) neighbours on the patch grid and 0 otherwise, and
 129 \bar{f}_i the mean activation across patches. A feature is *evaluable* on an image if it activates on at least
 130 5 patches with activation variance above $\varepsilon=10^{-8}$, and *evaluable in aggregate* if it is evaluable on

131 at least 50 validation images. The per-feature score \bar{I}_i is the mean of $I(\mathbf{f}_i)$ over the feature’s valid
 132 images. The SAE-level statistic is the mean across evaluable features:

$$\bar{I}_{\text{SAE}} = \frac{1}{N_{\text{eval}}} \sum_{i \in \mathcal{E}} \bar{I}_i, \quad (3)$$

133 where \mathcal{E} is the set of evaluable features and $N_{\text{eval}} = |\mathcal{E}|$. We report \bar{I}_{SAE} as the primary M1 score,
 134 alongside the median, interquartile range, and $N_{\text{eval}}/N_{\text{dict}}$ as the fraction of evaluable features.
 135 Higher values indicate more spatially coherent features.

136 **Sanity check against the spatial-randomness null.** Under the null hypothesis of no spatial
 137 autocorrelation, Moran’s I has analytic expectation $\mathbb{E}[I] = -1/(P-1)$ [Moran, 1950]. All 60 SAEs
 138 reject this null at $|z| > 2$; we therefore use z only as a sanity check and report mean Moran’s I , rather
 139 than z , as the M1 effect-size statistic because N_{eval} mechanically inflates z through the standard error.
 140 The exact z -statistic is given in Appendix F.

141 4.2 Reconstruction: FVU (M2) and Downstream Preservation (M3)

142 **Motivation.** Reconstruction quality is the most commonly reported SAE metric, almost always as
 143 Fraction of Variance Unexplained (FVU). The implicit assumption is that low FVU implies preserva-
 144 tion of task-relevant information. We argue this assumption fails empirically across backbones (§5)
 145 and therefore report FVU as M2 alongside a task-conditional preservation metric M3.

146 **Definition (M2: FVU).** Let $\mathbf{a} \in \mathbb{R}^D$ denote a backbone activation vector and $\hat{\mathbf{a}}$ its SAE reconstruc-
 147 tion. FVU is the fraction of activation variance unexplained:

$$\text{FVU} = \frac{\mathbb{E}_{\mathbf{a}} \left[\|\mathbf{a} - \hat{\mathbf{a}}\|_2^2 \right]}{\mathbb{E}_{\mathbf{a}} \left[\|\mathbf{a} - \bar{\mathbf{a}}\|_2^2 \right]}, \quad (4)$$

148 with $\bar{\mathbf{a}}$ the mean activation over the evaluation set. We compute FVU on the ImageNet-1K validation
 149 split using activations from the SAE training layer (§3). Lower is better.

150 **Definition (M3: Downstream Preservation).** Let ϕ denote a linear classifier trained on raw
 151 backbone activations, and let a and \hat{a} denote raw and SAE-reconstructed activations from a held-out
 152 evaluation split. We define

$$M3 = \frac{\text{acc}(\phi(\hat{a}))}{\text{acc}(\phi(a))}.$$

153 Unlike M2, which penalizes all activation-space reconstruction error, M3 measures the fraction
 154 of downstream-relevant information preserved after SAE reconstruction. All probes use matched
 155 architectures, splits, and hyperparameters; details are in Appendix A.

156 4.3 Concept Detection: Sparse Probing (M4), Monosemanticity (M5), Cross-Domain (M6)

157 **Motivation.** Spatial coherence (M1) and reconstruction (M2, M3) say nothing about whether SAE
 158 features are useful units for reading concepts out of the representation. M4 asks whether features
 159 support concept readout under sparsity constraints, M5 asks whether individual features correspond
 160 to consistent semantic content, and M6 asks whether these properties hold beyond the SAE’s training
 161 distribution. Concept detection is the dimension on which prior vision SAE work has concentrated
 162 [Pach et al., 2025] and the dimension along which H2 and H4 are tested (§4.5).

163 **Sparse Probing (M4).** Given an SAE with feature dimension $F \in \{8D, 16D, 32D\}$, let $c \in$
 164 \mathbb{R}^F denote the pooled SAE code for an image. For each of the 1000 ImageNet-1K classes, we
 165 rank SAE features by a univariate class-discrimination score on the probe-training split, train a
 166 logistic-regression classifier on the top- k ranked features, and evaluate AUC on held-out examples.
 167 M4 is reported as mean AUC at $k = 128$; the full sparsity–accuracy curve appears in Figure 6.
 168 Implementation details, including the ANOVA F-statistic ranking, split construction, and probe
 169 hyperparameters, are reported in Appendix A.

170 A non-sparse linear probe trained on all F features can recover most of the information in raw
 171 activations regardless of how features are organized, since the probe can find a linear combination that
 172 approximates the raw representation. The interesting question is whether individual SAE features,
 173 used in small numbers, support concept readout. The shape of the sparse-probing curve across feature
 174 budgets is therefore as informative as any k value, and we use $k = 128$ as the main summary point.

175 **Monosemanticity (M5).** For each alive SAE feature i , we identify its top- N maximally activat-
 176 ing validation images ($N=16$), embed them with a *separate* backbone, and define the feature’s
 177 monosemanticity score as the mean pairwise cosine similarity of the N embeddings. Aggregated to a
 178 per-SAE score, raw monosemanticity is normalized against a random baseline:

$$M5 = \frac{s_{\text{obs}} - s_{\text{baseline}}}{1 - s_{\text{baseline}}}, \quad (5)$$

179 where s_{baseline} is the mean pairwise cosine similarity of N randomly sampled validation images
 180 under the cross-model embedding. $M5=0$ means the feature’s top images are no more similar than
 181 random; $M5=1$ means maximally similar. To avoid circularity from scoring features in the same
 182 embedding space used to train the SAE, each SAE is evaluated with a different backbone, and the
 183 score is normalized by a random-image baseline; details are in Appendix K.

184 **Cross-Domain Generalization (M6).** SAEs are trained on ImageNet-1K activations and may
 185 overfit to ImageNet-distribution structure. M6 evaluates whether SAE representations preserve
 186 task-relevant information under distribution shift. We evaluate three non-ImageNet classification
 187 datasets: iNaturalist, CUB-200, and DTD, covering fine-grained species recognition, fine-grained
 188 bird classification, and texture recognition. For each dataset, matched linear probes are trained on raw
 189 backbone features, SAE-reconstructed features, and SAE code features. We report preservation as the
 190 reconstructed-feature probe accuracy divided by the raw-feature probe accuracy, and also evaluate
 191 sparse SAE-code probing at $k = 128$. This ratio factors out the raw OOD capacity of the backbone
 192 and isolates SAE-induced information loss. Full probe details are given in Appendix L. Full raw,
 193 reconstructed, and preserved accuracies appear in Appendix M.

194 4.4 Disentanglement: Probe-Based Feature Absorption Rate (M7)

195 **Motivation.** A well-disentangled SAE should expose class-relevant concepts through small numbers
 196 of individual features rather than requiring several related features to recover the same distinction.
 197 Feature absorption [Chanin et al., 2024] describes the failure mode in which information that could
 198 ideally be represented by a single feature is instead distributed across related features. We adapt this
 199 idea to vision through a probe-based criterion: a class distinction is treated as absorbed when the top
 200 individual feature is insufficient, but a small set of top-ranked features recovers the distinction.

201 **Definition.** For each ImageNet hierarchy sibling group with at least three classes, we construct
 202 class-versus-sibling probing tasks. SAE codes are pooled to image-level features, and features are
 203 ranked by a univariate class-discrimination score on the probe-training split. For each task, we train
 204 sparse logistic-regression probes using the top $k = 1$ and top $k = 4$ ranked features and compute
 205 their F1 scores. A task is counted as absorbed when

$$F1(k = 1) < \alpha \quad \text{and} \quad F1(k = 4) > \beta,$$

206 with $\alpha = 0.8$ and $\beta = 0.5$ by default. M7 is the fraction of eligible class-versus-sibling tasks
 207 satisfying this criterion. Lower values indicate fewer cases where a class distinction is unrecoverable
 208 from the top individual feature but recoverable from a small set of related features. Sensitivity to α
 209 and β is reported in Appendix N.

210 4.5 Hypotheses

211 The seven metrics measure properties internal to the SAE; whether those properties predict down-
 212 stream task utility is a separate empirical question. We pre-registered four hypotheses linking
 213 SAE-internal metrics to downstream tasks before computing cross-task correlations, each tested
 214 cross-backbone with $n = 5$ by aggregating the SAE-internal metric across the 12 configurations
 215 per backbone and correlating against a per-backbone downstream number. Spearman’s ρ is the test
 216 statistic, with confirmation thresholds fixed in advance: **H1** $M1 \rightarrow \text{ADE20K segmentation mIoU}$

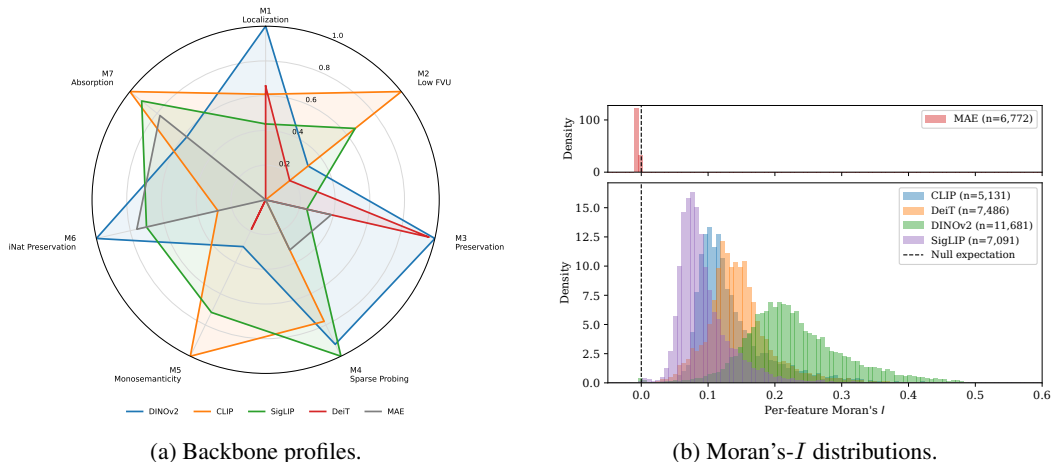


Figure 1: **Cross-backbone SAE behavior.** (a) Normalized ViSAEBench scores show distinct metric profiles across backbones. (b) Per-feature Moran's I separates MAE, which concentrates near the spatial-random null, from the other four spatially structured backbones.

217 ($\rho \geq 0.7$); **H2** normalized M5 \rightarrow CUB-200 fine-grained accuracy ($\rho \geq 0.6$); **H3** M3 \rightarrow image-image
 218 retrieval R@1 on a held-out ImageNet subset ($\rho \geq 0.6$); **H4** M4 (ImageNet) \rightarrow sparse probing AUC
 219 on iNaturalist at $k=128$ ($\rho \geq 0.7$), with CUB-200 and DTD as additional cross-domain checks.
 220 The mechanisms are: spatially coherent features carry the structure a segmentation head exploits
 221 (H1); within-feature consistency under an independent embedding should track semantically coherent
 222 visual content (H2); retrieval is a near-direct test of representation similarity (H3); SAE features that
 223 are concept-discriminative on ImageNet should remain so on related natural-image tasks (H4). All
 224 four results are reported in §5 regardless of outcome. No single scalar adequately captures vision
 225 SAE quality across the four dimensions, and the framework's value lies in characterizing trade-offs
 226 rather than producing a leaderboard.

227 5 Experiments and Findings

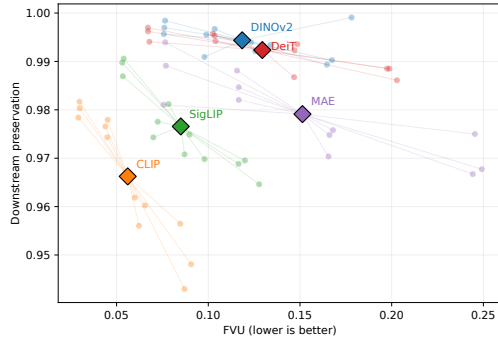
228 We trained 60 BatchTopK SAEs on identical 1.28M-image ImageNet-1K activation caches (5
 229 backbones, 3 expansion factors, 4 sparsity levels) and evaluated each along the seven metrics
 230 defined in §4. Figures 1a through 3 summarize the empirical findings.

231 **Cross-backbone profile** Figure 1a shows normalized scores across the seven metrics for each back-
 232 bone. The five backbones produce qualitatively different SAE feature profiles: DINOv2 dominates
 233 the spatial dimension (M1) and ranks first or second on concept-detection metrics; CLIP and SigLIP
 234 cluster together on monosemanticity and probing; DeiT trails on concept detection but achieves the
 235 highest reconstruction preservation; MAE ranks last on five of seven metrics. No backbone wins
 236 across all dimensions. The profiles are consistent with the view that pretraining objective shapes
 237 which kinds of structure are SAE-decomposable, rather than producing a uniform “better” or “worse”
 238 backbone for interpretability.

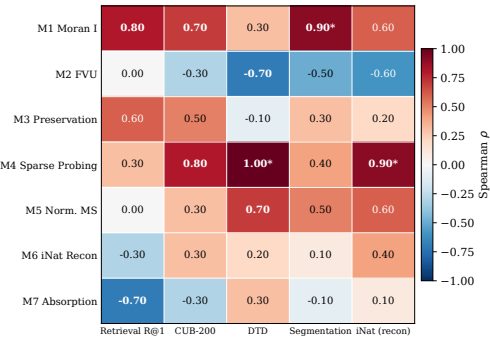
239 5.1 Finding 1: Categorical separation in spatial coherence

240 Four of the five backbones produce SAE features with strongly non-random spatial structure. M1
 241 (mean Moran's I) values are 0.232 (DINOv2), 0.149 (DeiT), 0.139 (CLIP), and 0.099 (SigLIP), all
 242 rejecting the spatial-randomness null $\mathbb{E}[I] = -1/(P-1)$ at $|z| > 150$. MAE produces SAE features
 243 whose spatial activation patterns are at the spatial-randomness null in effect size: $M1 = -0.005$,
 244 matching the analytic null expectation $\mathbb{E}[I] = -1/(P-1) = -0.005128$ (for $P=196$) to four
 245 decimals, across all 12 MAE configurations.

246 The interquartile range of MAE's per-feature distribution is 0.0008, two orders of magnitude tighter
 247 than DINOv2's 0.091, ruling out the possibility that flatness reflects cancellation between positively



(a) FVU preservation scatter.



(b) Metric-task correlation heatmap.

Figure 2: **Metric dissociations.** (a) FVU does not reliably predict downstream preservation across backbones. (b) Cross-backbone metric–task correlations show that different ViSAEBench metrics capture distinct downstream-relevant properties.

248 and negatively localized features. Figure 1b shows the per-feature distributions against the null. A
 249 diagnostic on raw pre-SAE backbone activations (Appendix G) confirms that MAE’s spatial flatness
 250 originates in the backbone, not the SAE. Raw MAE activations sit on the null at $\bar{I} = -0.005$. Raw
 251 activations from the other four backbones reject the null at $z > 150$, with mean Moran’s I between
 252 0.16 and 0.38. SAE decomposition preserves a substantial fraction of the raw spatial coherence
 253 (mean ratio 0.64 across the four non-null backbones), so the categorical gap exists in the source
 254 representation and the SAE inherits it. Figure 5 shows representative activation maps that visualize
 255 this gap directly.

256 This pattern is consistent with H1: the cross-backbone Spearman correlation between M1 and
 257 ADE20K segmentation mIoU is $\rho = 0.90$ ($p = 0.037$, $n = 5$), above the pre-registered threshold
 258 $\rho \geq 0.7$

259 5.2 Finding 2: FVU and downstream preservation dissociate

260 Figure 2a plots per-configuration FVU against downstream preservation. The two metrics are not
 261 monotonically related across backbones. CLIP achieves the lowest mean FVU of any backbone
 262 (0.056) but a mean preservation ratio of 0.966, the second-lowest. DeiT achieves the highest mean
 263 preservation (0.992) at FVU = 0.129. DINOv2 reaches 0.994 preservation at FVU = 0.118. Within-
 264 backbone Spearman correlation between FVU and preservation is weak across all five backbones,
 265 and the cross-backbone rank order produced by FVU inverts the order produced by preservation for
 266 the highest-fidelity backbones.

267 The dissociation has a clear practical consequence: comparing vision SAE quality across backbones
 268 using FVU alone produces misleading conclusions. CLIP’s contrastive pretraining concentrates label-
 269 aligned variance on directions used by the contrastive head, and the BatchTopK objective is unaware
 270 of which directions those are. Optimizing FVU under this objective can reduce reconstruction error in
 271 directions the linear probe never reads while damaging directions it does. Within a single backbone,
 272 FVU remains a useful relative metric for comparing SAE configurations. Across backbones, only
 273 task-conditional metrics (M3 or M6) support quality claims.

274 5.3 Finding 3: Sparse probing transfers across natural-image domains (H4)

275 Cross-backbone Spearman correlation between M4 (ImageNet sparse probing AUC at $k = 128$)
 276 and iNaturalist sparse probing AUC at $k = 128$ is $\rho = 0.90$ ($p = 0.037$), clearing the H4 threshold.
 277 The same trend holds on DTD ($\rho = 1.00$) and CUB-200 ($\rho = 0.80$), with a weaker association to
 278 segmentation transfer ($\rho = 0.40$). M4 is the strongest single predictor of cross-task generalization
 279 in our suite (Figure 2b): backbones whose SAE features are concept-discriminative on ImageNet
 280 remain so on related natural-image tasks.

281 **5.4 Finding 4: Preservation predicts retrieval but below threshold (H3)**

282 Cross-backbone Spearman correlation between M3 and image-image retrieval R@1 is $\rho = 0.60$
283 ($p = 0.285$), meeting but not exceeding the H3 threshold given the small sample size. The directional
284 prediction holds, and Figure 2b shows M1 ($\rho = 0.80$) and M3 ($\rho = 0.60$) jointly tracking retrieval.
285 We treat H3 as supported but underpowered.

286 **5.5 Finding 5: Monosemanticity does not predict fine-grained classification**

287 Monosemanticity, as operationalized by top-activating-image consistency under an independent
288 embedding, does not predict fine-grained classification ability across backbones. Cross-backbone
289 Spearman correlation between normalized M5 and CUB-200 accuracy is $\rho = 0.30$ ($p = 0.624$,
290 $n = 5$), below the pre-registered threshold of 0.6, with the directional prediction also unsupported.

291 The result has a clean mechanism. Within-feature consistency is a property of a feature’s activation
292 set in isolation: how perceptually similar are the images on which feature i fires most strongly.
293 Between-class separability is a property of the joint geometry of features across classes: do the
294 features collectively distinguish bird species A from species B. A feature can fire consistently on
295 visually similar images that span multiple fine-grained categories, registering as highly monosemantic
296 while contributing nothing to fine-grained discrimination. The two properties are formally orthogonal,
297 and our data show they are also empirically dissociated.

298 We decompose M5 into within-class and between-class components by computing mean pairwise
299 cosine similarity within the same ImageNet validation class versus across ImageNet classes. High-M5
300 features often retain substantial between-class similarity, confirming that M5 captures visual consis-
301 tency but not necessarily fine-grained class separability (Appendix K). Figure 5 shows representative
302 high-M5 features whose top activating images are perceptually consistent (the bird-on-branch case)
303 but straddle multiple fine-grained classes.

304 This finding has implications for SAE evaluation practice. Monosemanticity is a central quality
305 criterion in language SAE work [Bricken et al., 2023, Templeton et al., 2024]; the dissociation
306 we observe in vision raises the question of whether the metric tracks downstream task utility in
307 language either. We recommend reporting M5 jointly with concept-detection metrics that test between-
308 feature distinctness, and we suggest that future monosemanticity metrics incorporate between-class
309 separability, for example through a discriminability-weighted consistency score, rather than treating
310 within-feature consistency as a sufficient quality signal.

311 **5.6 Finding 6: DeiT’s absorption paradox**

312 DeiT achieves the lowest probe-based absorption rate under the default thresholds ($\alpha = 0.8$, $\beta = 0.5$),
313 with $M7 = 0.298$ compared to 0.415 (MAE), 0.423 (DINOv2), 0.494 (CLIP), and 0.547 (SigLIP),
314 yet produces the second-worst sparse probing AUC, so low absorption does not imply high concept-
315 detection performance. DeiT remains the lowest-absorption backbone in 7 of 9 settings under a
316 sweep of $\alpha \in \{0.7, 0.8, 0.9\}$ and $\beta \in \{0.3, 0.5, 0.7\}$, breaking only in the strictest corner ($\alpha =$
317 0.9 , $\beta \geq 0.5$) where MAE and DeiT converge; the paradox is a stable directional pattern rather than
318 an ordering invariant (Appendix N). Recent work on decoder orthogonality constraints [Korznikov
319 et al., 2025] suggests absorption-prone failure modes may be addressable architecturally, motivating
320 future investigation into whether orthogonal decoders reconcile DeiT’s behavior.

321 **5.7 Variance decomposition**

322 A natural concern with cross-architecture SAE benchmarking is whether the metric differences reflect
323 SAE properties or are predictable from the pretrained backbone. We address this with a two-way
324 ANOVA on each metric, with backbone (5 levels) and hyperparameter configuration (expansion factor
325 $\times K$, 12 levels) as factors.

326 Table 3 (Appendix B) reports η^2 , the proportion of total sum-of-squares attributable to each factor.

327 As shown in Fig. 3, backbone explains over 90% of total variance on M1, M5, and M7, and roughly
328 three-quarters on M4 and M6. Differences in expansion factor and K account for almost none of
329 the cross-configuration spread on spatial coherence, monosemanticity, or absorption. M2 inverts the
330 pattern: backbone explains 39.2% of FVU variance while hyperparameters explain 52.3%. M3 sits

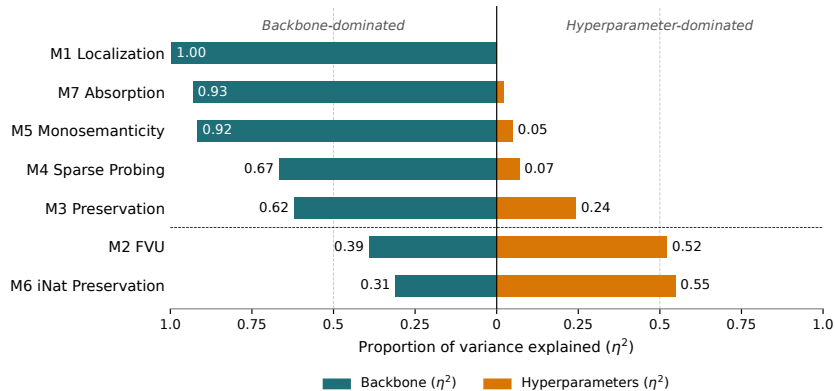


Figure 3: Variance decomposition of each metric across the 60 SAE configurations.

331 between, with backbone at 62.2% and hyperparameters at 24.3%. Restricting the ANOVA to the
 332 four non-MAE backbones yields $\eta_{\text{backbone}}^2 = 0.987$ for M1, confirming that the headline result is not
 333 driven by MAE’s null-indistinguishable cluster.

334 5.8 Hypothesis summary

335 Of the four pre-registered hypotheses, three are confirmed (H1, H3, H4). H1 and H4 clear the
 336 pre-registered threshold at $\rho \geq 0.7$. H3 meets but does not exceed its threshold. The fourth, H2, is
 337 not supported: monosemanticity, as commonly defined, does not track fine-grained classification
 338 ability across backbones, with implications for SAE evaluation practice in both vision and language.
 339 Figure 2b reports the full 7×5 correlation table across all metric–task pairs, including the 31
 340 non-pre-registered combinations.

341 6 Limitations and Open Questions

342 ViSAEBench evaluates five ViT-B backbones at a single layer with a single SAE family (BatchTopK).
 343 The $n = 5$ backbone count limits the statistical power of cross-backbone correlations, and we treat
 344 single-task Spearman p-values as supporting rather than confirmatory evidence. SAEs are trained
 345 on ImageNet-1K activations, so features specific to non-natural-image domains (medical imaging,
 346 satellite imagery, synthetic data) may be underfit, and other SAE architectures (JumpReLU, Gated
 347 SAE, Matryoshka) may exhibit different cross-backbone profiles. M1 measures whether features
 348 activate in spatially contiguous regions, not whether those regions are semantically meaningful
 349 (Figure 5). Methodologically, the variance decomposition (§5.7) implies that single-backbone
 350 evaluations of spatial coherence, monosemanticity, and absorption measure backbone properties more
 351 than SAE properties; we recommend backbone standardization as a prerequisite for SAE-method
 352 comparison on these dimensions. FVU is the exception, with hyperparameter choice as its primary
 353 driver, making cross-backbone reporting comparatively safer there.

354 7 Conclusion

355 We introduced ViSAEBench, an evaluation suite spanning seven metrics across four interpretability
 356 dimensions, and applied it to 60 BatchTopK SAEs across five ViT-B backbones. The cross-backbone
 357 study surfaces three dissociations that constrain how vision SAEs should be evaluated: MAE-derived
 358 SAEs are statistically indistinguishable from a spatially-random null while the other four backbones
 359 are not; FVU and downstream preservation diverge across backbones; and monosemanticity does
 360 not predict fine-grained classification, despite its central role in language SAE work. Three of four
 361 pre-registered hypotheses are confirmed, and the fourth yields a substantive negative result with
 362 implications for SAE evaluation in both vision and language.

363 **References**

- 364 Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner,
365 Cem Anil, Carson Denison, Amanda Askeell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer,
366 Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean,
367 Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity:
368 Decomposing language models with dictionary learning. [https://transformer-circuits.pub/2023/
369 monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html), 2023.
- 370 Bart Bussmann, Patrick Leask, and Neel Nanda. Batchtopk sparse autoencoders. *arXiv preprint*
371 *arXiv:2412.06410*, 2024.
- 372 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin.
373 Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International*
374 *Conference on Computer Vision (ICCV)*, pages 9650–9660, October 2021.
- 375 David Chanin, James Wilken-Smith, Tomáš Dulka, Hardik Bhatnagar, Satvik Golechha, and Joseph Bloom.
376 A is for absorption: Studying feature splitting and absorption in sparse autoencoders. *arXiv preprint*
377 *arXiv:2409.14507*, 2024.
- 378 Gerasimos Chatzoudis, Konstantinos D. Polyzos, Zhuowei Li, Difei Gu, Gemma E. Moran, Hao Wang, and
379 Dimitris N. Metaxas. Can cross-layer transcoders replace vision transformer activations? an interpretable
380 perspective on vision, 2026. URL <https://arxiv.org/abs/2604.13304>.
- 381 Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers.
382 In *The Twelfth International Conference on Learning Representations*, 2024. URL [https://openreview.
383 net/forum?id=2dn03LLiJ1](https://openreview.net/forum?id=2dn03LLiJ1).
- 384 Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike,
385 and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint arXiv:2406.04093*, 2024.
- 386 Difei Gu, Yunhe Gao, Gerasimos Chatzoudis, Zihan Dong, Guoning Zhang, Bangwei Guo, Yang Zhou, Mu Zhou,
387 and Dimitris Metaxas. Lucid-sae: Learning unified vision-language sparse codes for interpretable concept
388 discovery, 2026. URL <https://arxiv.org/abs/2602.07311>.
- 389 Sangyu Han, Yearim Kim, and Nojun Kwak. Causal interpretation of sparse autoencoder features in vision,
390 2025. URL <https://arxiv.org/abs/2509.00749>.
- 391 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are
392 scalable vision learners. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
393 2022. URL <https://arxiv.org/abs/2111.06377>.
- 394 Sonia Joseph, Praneet Suresh, Lorenz Hufe, Edward Stevinson, Robert Graham, Yash Vadi, Danilo Bzdok, Se-
395 bastian Lapuschkin, Lee Sharkey, and Blake Aaron Richards. Prisma: An open source toolkit for mechanistic
396 interpretability in vision and video, 2025. URL <https://arxiv.org/abs/2504.19475>.
- 397 Adam Karvonen, Can Rager, Johnny Lin, Curt Tigges, Joseph Bloom, David Chanin, Yeu-Tong Lau, Eoin Farrell,
398 Callum McDougall, Kola Ayonrinde, Demian Till, Matthew Wearden, Arthur Conmy, Samuel Marks, and Neel
399 Nanda. Saebench: A comprehensive benchmark for sparse autoencoders in language model interpretability,
400 2025. URL <https://arxiv.org/abs/2503.09532>.
- 401 Anton Korznikov, Andrey Galichin, Alexey Dontsov, Oleg Rogov, Elena Tutubalina, and Ivan Oseledets. Ortsae:
402 Orthogonal sparse autoencoders uncover atomic features, 2025. URL [https://arxiv.org/abs/2509.
403 22033](https://arxiv.org/abs/2509.22033).
- 404 Tom Lieberum, Senthoran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma,
405 János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse autoencoders
406 everywhere all at once on gemma 2, 2024. URL <https://arxiv.org/abs/2408.05147>.
- 407 Alireza Makhzani and Brendan Frey. A winner-take-all method for training sparse convolutional autoencoders.
408 In *NIPS deep learning workshop*, 2014.
- 409 P. A. P. Moran. Notes on continuous stochastic phenomena. *Biometrika*, 37(1/2):17–23, 1950. ISSN 00063444,
410 14643510. URL <http://www.jstor.org/stable/2332142>.
- 411 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez,
412 Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba,
413 Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve,
414 Hu Xu, Hervé Jégou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning
415 robust visual features without supervision, 2024. URL <https://arxiv.org/abs/2304.07193>.

- 416 Mateusz Pach, Shyamgopal Karthik, Quentin Bouniot, Serge Belongie, and Zeynep Akata. Sparse autoencoders
417 learn monosemantic features in vision-language models, 2025. URL [https://arxiv.org/abs/2504.](https://arxiv.org/abs/2504.02821)
418 02821.
- 419 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
420 Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable
421 visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- 422 Senthoooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János Kramár,
423 and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse autoencoders, 2024.
424 URL <https://arxiv.org/abs/2407.14435>.
- 425 Alexander Rubinstein, Ameya Prabhu, Matthias Bethge, and Seong Joon Oh. Are we done with object-centric
426 learning?, 2025. URL <https://arxiv.org/abs/2504.07092>.
- 427 Lee Sharkey, Bilal Chughtai, Joshua Batson, Jack Lindsey, Jeff Wu, Lucius Bushnaq, Nicholas Goldowsky-Dill,
428 Stefan Heimersheim, Alejandro Ortega, Joseph Bloom, Stella Biderman, Adria Garriga-Alonso, Arthur
429 Conmy, Neel Nanda, Jessica Rumbelow, Martin Wattenberg, Nandi Schoots, Joseph Miller, Eric J. Michaud,
430 Stephen Casper, Max Tegmark, William Saunders, David Bau, Eric Todd, Atticus Geiger, Mor Geva, Jesse
431 Hoogland, Daniel Murfet, and Tom McGrath. Open problems in mechanistic interpretability, 2025. URL
432 <https://arxiv.org/abs/2501.16496>.
- 433 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce,
434 Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L. Turner, Callum McDougall,
435 Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn,
436 Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from
437 claude 3 sonnet. [https://transformer-circuits.pub/2024/scaling-monosemanticity/index.](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html)
438 [html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html), 2024. Accessed: 2026-05-03.
- 439 Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou.
440 Training data-efficient image transformers & distillation through attention, 2021. URL [https://arxiv.](https://arxiv.org/abs/2012.12877)
441 [org/abs/2012.12877](https://arxiv.org/abs/2012.12877).
- 442 Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruysen, Carlos Riquelme, Mario Lucic, Josip
443 Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, Lucas Beyer, Olivier Bachem, Michael
444 Tschannen, Marcin Michalski, Olivier Bousquet, Sylvain Gelly, and Neil Houlsby. A large-scale study of
445 representation learning with the visual task adaptation benchmark, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1910.04867)
446 [1910.04867](https://arxiv.org/abs/1910.04867).
- 447 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language-image
448 pre-training. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. URL [https:](https://arxiv.org/abs/2303.15343)
449 [//arxiv.org/abs/2303.15343](https://arxiv.org/abs/2303.15343).
- 450 Theodore Zhengde Zhao, Sid Kiblawi, Jianwei Yang, Naoto Usuyama, Reuben Tan, Noel C Codella, Tristan Nau-
451 mann, Hoifung Poon, and Mu Wei. Learning sparse visual representations via spatial-semantic factorization,
452 2026. URL <https://arxiv.org/abs/2602.01905>.

453 **NeurIPS Paper Checklist**

454 **1. Claims**

455 Question: Do the main claims made in the abstract and introduction accurately reflect the
456 paper’s contributions and scope?

457 Answer: [\[Yes\]](#)

458 Justification: The abstract and introduction state three concrete claims: (1) ViSAEBench
459 provides seven metrics across four interpretability dimensions including a novel spatial
460 coherence metric; (2) backbone explains over 90% of variance on three of seven metrics and
461 over 60% on four of seven, with FVU being the sole hyperparameter-dominated exception;
462 (3) MAE-derived SAEs are categorically separated from the other four backbones in spatial
463 coherence at $z > 150$. All three claims are supported by experiments in Section 5 and the
464 variance decomposition in Section 5.7 and Table 3. The framing as a benchmarking suite
465 rather than a leaderboard is reflected in the explicit “no composite score” position in Section
466 4.5.

467 Guidelines:

- 468 • The answer [N/A] means that the abstract and introduction do not include the claims
469 made in the paper.
- 470 • The abstract and/or introduction should clearly state the claims made, including the
471 contributions made in the paper and important assumptions and limitations. A [No] or
472 [N/A] answer to this question will not be perceived well by the reviewers.
- 473 • The claims made should match theoretical and experimental results, and reflect how
474 much the results can be expected to generalize to other settings.
- 475 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
476 are not attained by the paper.

477 2. Limitations

478 Question: Does the paper discuss the limitations of the work performed by the authors?

479 Answer: [Yes]

480 Justification: Section 6 discusses limitations explicitly: (1) scope is restricted to five ViT-B
481 backbones at a single layer with a single SAE family (BatchTopK), with n=5 limiting
482 cross-backbone correlation power; (2) SAEs are trained on ImageNet-1K and may underfit
483 non-natural-image domains such as medical imaging or satellite imagery; (3) other SAE
484 architectures including JumpReLU, Gated, and Matryoshka are not benchmarked; (4) the
485 MAE mechanism is reported empirically without direct manipulation of the pretraining
486 objective; (5) M1 measures spatial contiguity but not semantic meaningfulness and must be
487 read jointly with M4. The paper also acknowledges single-task Spearman pp p-values are
488 treated as supporting rather than confirmatory given n=5n=5 n=5.

489 Guidelines:

- 490 • The answer [N/A] means that the paper has no limitation while the answer [No] means
491 that the paper has limitations, but those are not discussed in the paper.
- 492 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 493 • The paper should point out any strong assumptions and how robust the results are to
494 violations of these assumptions (e.g., independence assumptions, noiseless settings,
495 model well-specification, asymptotic approximations only holding locally). The authors
496 should reflect on how these assumptions might be violated in practice and what the
497 implications would be.
- 498 • The authors should reflect on the scope of the claims made, e.g., if the approach was
499 only tested on a few datasets or with a few runs. In general, empirical results often
500 depend on implicit assumptions, which should be articulated.
- 501 • The authors should reflect on the factors that influence the performance of the approach.
502 For example, a facial recognition algorithm may perform poorly when image resolution
503 is low or images are taken in low lighting. Or a speech-to-text system might not be
504 used reliably to provide closed captions for online lectures because it fails to handle
505 technical jargon.
- 506 • The authors should discuss the computational efficiency of the proposed algorithms
507 and how they scale with dataset size.
- 508 • If applicable, the authors should discuss possible limitations of their approach to
509 address problems of privacy and fairness.
- 510 • While the authors might fear that complete honesty about limitations might be used by
511 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
512 limitations that aren’t acknowledged in the paper. The authors should use their best
513 judgment and recognize that individual actions in favor of transparency play an impor-
514 tant role in developing norms that preserve the integrity of the community. Reviewers
515 will be specifically instructed to not penalize honesty concerning limitations.

516 3. Theory assumptions and proofs

517 Question: For each theoretical result, does the paper provide the full set of assumptions and
518 a complete (and correct) proof?

519 Answer: [N/A]

520 Justification: The paper does not contain theoretical results requiring proofs. M1 uses
521 Moran’s I, which is a standard spatial statistics quantity with a well-known null distribution;
522 the formula is stated and cited but no novel theorems are claimed.

523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Section 3.3 specifies the SAE training protocol: 12 configurations per backbone sweeping expansion factors {8, 16, 32} and $K \in \{64, 128, 192, 256\}$ on 1.28M ImageNet-1K activations from block 11 of each backbone, with identical optimizer, learning rate, and step counts across backbones. Backbones are listed in Table 1 with sources. Section 4 defines all seven metrics formally, with thresholds, parameter values ($N = 16$ for M5, $\alpha = 0.8$ and $\beta = 0.5$ for M7, $\varepsilon = 10^{-8}$ for M1 evaluability), and evaluation datasets. Full hyperparameters and layer-sensitivity analysis are deferred to Appendix A and C. All 60 SAE checkpoints are released on HuggingFace under the visaebench organization.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

579 **5. Open access to data and code**

580 Question: Does the paper provide open access to the data and code, with sufficient instruc-
581 tions to faithfully reproduce the main experimental results, as described in supplemental
582 material?

583 Answer: [Yes]

584 Justification: All 60 SAE checkpoints will be released on HuggingFace under the visaebench
585 organization in safetensors format under Apache 2.0 license, with one repository per back-
586 bone containing the 12 configuration subfolders. The evaluation library is released will
587 be released as a pip-installable python package once the paper is accepted. ImageNet-1K,
588 iNaturalist 2021, CUB-200-2011, DTD, and ADE20K are all public datasets.

589 Guidelines:

- 590 • The answer [N/A] means that paper does not include experiments requiring code.
- 591 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
592 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 593 • While we encourage the release of code and data, we understand that this might not
594 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
595 including code, unless this is central to the contribution (e.g., for a new open-source
596 benchmark).
- 597 • The instructions should contain the exact command and environment needed to run to
598 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 599 • The authors should provide instructions on data access and preparation, including how
600 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 601 • The authors should provide scripts to reproduce all experimental results for the new
602 proposed method and baselines. If only a subset of experiments are reproducible, they
603 should state which ones are omitted from the script and why.
- 604 • At submission time, to preserve anonymity, the authors should release anonymized
605 versions (if applicable).
- 606 • Providing as much information as possible in supplemental material (appended to the
607 paper) is recommended, but including URLs to data and code is permitted.
- 608

609 **6. Experimental setting/details**

610 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
611 rameters, how they were chosen, type of optimizer) necessary to understand the results?

612 Answer: [Yes]

613 Justification: Section 3.3 specifies the SAE sweep grid (expansion factors, K values, 1.28M
614 training activations, block 11 tap point, identical optimizer settings across backbones).
615 Section 4 specifies metric-level details: M1 evaluability thresholds (5 patches, variance
616 above 10^{-8} , 50 evaluable images), M4 selection rule (top- k by mutual information at
617 $k \in \{32, 128, 512\}$), M5 cross-model embedding pairing with $N = 16$ maximally activating
618 images, M6 datasets (iNat, CUB, DTD, ADE20K) with EuroSAT exclusion documented,
619 and M7 thresholds $\alpha = 0.8$, $\beta = 0.5$. Train/test splits follow the standard ImageNet-1K
620 validation split for training-distribution evaluation. Full optimizer hyperparameters and
621 layer-sensitivity analyses are in Appendix A and C.

622 Guidelines:

- 623 • The answer [N/A] means that the paper does not include experiments.
- 624 • The experimental setting should be presented in the core of the paper to a level of detail
625 that is necessary to appreciate the results and make sense of them.
- 626 • The full details can be provided either with the code, in appendix, or as supplemental
627 material.

628 **7. Experiment statistical significance**

629 Question: Does the paper report error bars suitably and correctly defined or other appropriate
630 information about the statistical significance of the experiments?

631 Answer: [Yes]

632 Justification: The paper reports statistical significance in three ways. (1) M1 is itself a
633 z -statistic against a permutation null with the null expectation $\mathbb{E}[\bar{I}_{\text{null}}] = -1/(P - 1)$ stated
634 explicitly, and individual MAE configurations are flagged as null-indistinguishable. (2) The
635 four pre-registered hypotheses report Spearman ρ with p -values ($n = 5$): H1 $\rho = 0.90$,
636 $p = 0.037$; H2 $\rho = 0.30$, $p = 0.624$; H3 $\rho = 0.60$, $p = 0.285$; H4 $\rho = 0.90$, $p = 0.037$.
637 (3) The two-way ANOVA in Section 5.7 reports η^2 for backbone and hyperparameter factors
638 across all seven metrics in Table 3. The paper is also explicit that with $n = 5$ backbones,
639 single-task p -values are treated as supporting rather than confirmatory, with consistency
640 across multiple downstream tasks Figure 2b as primary evidence.

641 Guidelines:

- 642 • The answer [N/A] means that the paper does not include experiments.
- 643 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
644 intervals, or statistical significance tests, at least for the experiments that support the
645 main claims of the paper.
- 646 • The factors of variability that the error bars are capturing should be clearly stated (for
647 example, train/test split, initialization, random drawing of some parameter, or overall
648 run with given experimental conditions).
- 649 • The method for calculating the error bars should be explained (closed form formula,
650 call to a library function, bootstrap, etc.)
- 651 • The assumptions made should be given (e.g., Normally distributed errors).
- 652 • It should be clear whether the error bar is the standard deviation or the standard error
653 of the mean.
- 654 • It is OK to report 1-sigma error bars, but one should state it. The authors should
655 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
656 of Normality of errors is not verified.
- 657 • For asymmetric distributions, the authors should be careful not to show in tables or
658 figures symmetric error bars that would yield results that are out of range (e.g., negative
659 error rates).
- 660 • If error bars are reported in tables or plots, the authors should explain in the text how
661 they were calculated and reference the corresponding figures or tables in the text.

662 8. Experiments compute resources

663 Question: For each experiment, does the paper provide sufficient information on the com-
664 puter resources (type of compute workers, memory, time of execution) needed to reproduce
665 the experiments?

666 Answer: [Yes]

667 Justification: Compute used: NYU Torch H200 cluster (1 node, 8 GPU) and ai4ce-shannon
668 cluster, with prototyping on a local RTX 3090 (24GB). Slurm with Singularity containers
669 was used for full sweeps. The full sweep is 60 SAEs (5 backbones \times 12 configurations)
670 trained on 1.28M ImageNet-1K activations. Per-SAE training time, total GPU-hours, and
671 storage requirements for the cached activation tensors should be reported in Appendix A.

672 Guidelines:

- 673 • The answer [N/A] means that the paper does not include experiments.
- 674 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
675 or cloud provider, including relevant memory and storage.
- 676 • The paper should provide the amount of compute required for each of the individual
677 experimental runs as well as estimate the total compute.
- 678 • The paper should disclose whether the full research project required more compute
679 than the experiments reported in the paper (e.g., preliminary or failed experiments that
680 didn't make it into the paper).

681 9. Code of ethics

682 Question: Does the research conducted in the paper conform, in every respect, with the
683 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

684 Answer: [Yes]

685 Justification: The research conforms to the NeurIPS Code of Ethics. The work uses publicly
686 available pretrained backbones and standard public datasets (ImageNet-1K, iNaturalist 2021,
687 CUB-200-2011, DTD, ADE20K), involves no human subjects or crowdsourcing, and the
688 released artifacts are SAE checkpoints and an evaluation library, neither of which presents
689 elevated misuse risk relative to the underlying backbones.

690 Guidelines:

- 691 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
692 Ethics.
- 693 • If the authors answer [No], they should explain the special circumstances that require a
694 deviation from the Code of Ethics.
- 695 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
696 eration due to laws or regulations in their jurisdiction).

697 10. Broader impacts

698 Question: Does the paper discuss both potential positive societal impacts and negative
699 societal impacts of the work performed?

700 Answer: [Yes]

701 Justification: The work is interpretability infrastructure for vision models. Positive impacts:
702 standardized cross-backbone evaluation supports more rigorous mechanistic interpretability
703 research, helps identify representational pathologies (such as MAE’s null-indistinguishable
704 spatial structure), and reduces wasted compute from non-comparable single-paper SAE
705 evaluations. Negative impacts are indirect: more reliable interpretability tooling could be
706 applied to vision systems in surveillance or content moderation pipelines, but the contribution
707 itself is methodological infrastructure rather than a deployable system, and SAE evaluation
708 does not enable capabilities beyond those already present in the underlying pretrained
709 backbones.

710 Guidelines:

- 711 • The answer [N/A] means that there is no societal impact of the work performed.
- 712 • If the authors answer [N/A] or [No], they should explain why their work has no societal
713 impact or why the paper does not address societal impact.
- 714 • Examples of negative societal impacts include potential malicious or unintended uses
715 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
716 (e.g., deployment of technologies that could make decisions that unfairly impact specific
717 groups), privacy considerations, and security considerations.
- 718 • The conference expects that many papers will be foundational research and not tied
719 to particular applications, let alone deployments. However, if there is a direct path to
720 any negative applications, the authors should point it out. For example, it is legitimate
721 to point out that an improvement in the quality of generative models could be used to
722 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
723 that a generic algorithm for optimizing neural networks could enable people to train
724 models that generate Deepfakes faster.
- 725 • The authors should consider possible harms that could arise when the technology is
726 being used as intended and functioning correctly, harms that could arise when the
727 technology is being used as intended but gives incorrect results, and harms following
728 from (intentional or unintentional) misuse of the technology.
- 729 • If there are negative societal impacts, the authors could also discuss possible mitigation
730 strategies (e.g., gated release of models, providing defenses in addition to attacks,
731 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
732 feedback over time, improving the efficiency and accessibility of ML).

733 11. Safeguards

734 Question: Does the paper describe safeguards that have been put in place for responsible
735 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
736 image generators, or scraped datasets)?

737 Answer: [N/A]

738 Justification: The released artifacts are SAE checkpoints trained on activations from publicly
739 released pretrained backbones, plus an evaluation library.

740 Guidelines:

- 741 • The answer [N/A] means that the paper poses no such risks.
- 742 • Released models that have a high risk for misuse or dual-use should be released with
743 necessary safeguards to allow for controlled use of the model, for example by requiring
744 that users adhere to usage guidelines or restrictions to access the model or implementing
745 safety filters.
- 746 • Datasets that have been scraped from the Internet could pose safety risks. The authors
747 should describe how they avoided releasing unsafe images.
- 748 • We recognize that providing effective safeguards is challenging, and many papers do
749 not require this, but we encourage authors to take this into account and make a best
750 faith effort.

751 12. Licenses for existing assets

752 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
753 the paper, properly credited and are the license and terms of use explicitly mentioned and
754 properly respected?

755 Answer: [Yes]

756 Justification: All five backbones are publicly released under their original licenses and
757 properly cited (Table 1): DINOv2 (Oquab et al., 2024, CC-BY-NC 4.0), MAE (He et al.,
758 2022, CC-BY-NC 4.0), CLIP (Radford et al., 2021, OpenAI License), SigLIP (Zhai et al.,
759 2023, Apache 2.0), DeiT (Touvron et al., 2021, Apache 2.0). Datasets used are ImageNet-1K
760 (custom research license), iNaturalist 2021 (CC-BY-NC 4.0), CUB-200-2011 (CC0 1.0),
761 DTD (CC-BY 4.0), and ADE20K (custom research license), all publicly available with their
762 respective licenses. The overcomplete library BatchTopK SAE (Apache 2.0) and Prisma
763 library (Joseph et al., 2025, MIT) are cited.

764 Guidelines:

- 765 • The answer [N/A] means that the paper does not use existing assets.
- 766 • The authors should cite the original paper that produced the code package or dataset.
- 767 • The authors should state which version of the asset is used and, if possible, include a
768 URL.
- 769 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 770 • For scraped data from a particular source (e.g., website), the copyright and terms of
771 service of that source should be provided.
- 772 • If assets are released, the license, copyright information, and terms of use in the
773 package should be provided. For popular datasets, `paperswithcode.com/datasets`
774 has curated licenses for some datasets. Their licensing guide can help determine the
775 license of a dataset.
- 776 • For existing datasets that are re-packaged, both the original license and the license of
777 the derived asset (if it has changed) should be provided.
- 778 • If this information is not available online, the authors are encouraged to reach out to
779 the asset’s creators.

780 13. New assets

781 Question: Are new assets introduced in the paper well documented and is the documentation
782 provided alongside the assets?

783 Answer: [Yes]

784 Justification: Two new assets are released: (1) 60 SAE checkpoints under the HuggingFace,
785 organized as five repositories (one per backbone) each containing the 12 configuration,
786 released under Apache 2.0 license; (2) a pip-installable evaluation library implementing
787 the seven metrics. Documentation of training hyperparameters, configuration naming, and
788 metric usage will be in the repository README and Appendix A and C.

789 Guidelines:

- 790 • The answer [N/A] means that the paper does not release new assets.

- 791
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - 792
 - 793
 - 794
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - 795
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
 - 796
 - 797

798 **14. Crowdsourcing and research with human subjects**

799 Question: For crowdsourcing experiments and research with human subjects, does the paper
800 include the full text of instructions given to participants and screenshots, if applicable, as
801 well as details about compensation (if any)?

802 Answer: [N/A]

803 Justification: The work does not involve crowdsourcing or human subjects. All evaluations
804 use publicly available pretrained backbones and standard public benchmark datasets.

805 Guidelines:

- 806 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
807 with human subjects.
- 808 • Including this information in the supplemental material is fine, but if the main contribu-
809 tion of the paper involves human subjects, then as much detail as possible should be
810 included in the main paper.
- 811 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
812 or other labor should be paid at least the minimum wage in the country of the data
813 collector.

814 **15. Institutional review board (IRB) approvals or equivalent for research with human
815 subjects**

816 Question: Does the paper describe potential risks incurred by study participants, whether
817 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
818 approvals (or an equivalent approval/review based on the requirements of your country or
819 institution) were obtained?

820 Answer: [N/A]

821 Justification: The work does not involve human subjects research and therefore does not
822 require IRB approval.

823 Guidelines:

- 824 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
825 with human subjects.
- 826 • Depending on the country in which research is conducted, IRB approval (or equivalent)
827 may be required for any human subjects research. If you obtained IRB approval, you
828 should clearly state this in the paper.
- 829 • We recognize that the procedures for this may vary significantly between institutions
830 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
831 guidelines for their institution.
- 832 • For initial submissions, do not include any information that would break anonymity (if
833 applicable), such as the institution conducting the review.

834 **16. Declaration of LLM usage**

835 Question: Does the paper describe the usage of LLMs if it is an important, original, or
836 non-standard component of the core methods in this research? Note that if the LLM is used
837 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
838 scientific rigor, or originality of the research, declaration is not required.

839 Answer: [N/A]

840 Justification:

841 Guidelines:

842 • The answer [N/A] means that the core method development in this research does not
843 involve LLMs as any important, original, or non-standard components.

844 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
845 be described.

846 Appendix

847 A Experimental Setup and Hyperparameters

848 We report the full hyperparameters used for SAE training and evaluation in Tables 1 and 2. Unless
849 otherwise stated, all values are held fixed across the five backbones and across all 60 SAE configura-
850 tions. The only swept SAE hyperparameters are the expansion factor and the BatchTopK sparsity
851 level K .

Table 1: SAE training hyperparameters used in the main sweep.

Hyperparameter	Value
SAE family	BatchTopK SAE
Backbones	CLIP ViT-B/16, DINOv2 ViT-B/14, SigLIP ViT-B/16, MAE ViT-B/16, DeiT ViT-B/16
Layer	11
Input dimension	768
Expansion factors	{8, 16, 32}
Dictionary sizes	{6144, 12288, 24576}
K values	{64, 128, 192, 256}
Training activations	Cached ImageNet-1K train activations from backbone’s layer-11 patch tokens
Activation normalization	Per-dimension mean subtraction and scalar standard-deviation division using cached activation statistics
Optimizer	Adam
Learning rate	1×10^{-3}
Batch size	8192 patch tokens per step
Epochs	1
Seed	42
Precision	FP16
Inference threshold	Fixed BatchTopK threshold estimated from validation activations and saved with each checkpoint

Table 2: Evaluation hyperparameters used for ViSAEBench metrics (M1–M7).

Metric	Hyperparameters / Protocol
M1 (Spatial coherence)	Grid inferred from patch count (14×14 for 196 patches and 16×16 for 256 patches); 8-neighbor connectivity; minimum active patches per image = 5; activation variance threshold $\epsilon = 10^{-8}$; minimum valid images per feature = 50; SAE encoding batch size = 512; image-group size = 64.
M2 (FVU)	Computed on cached ImageNet-1K validation activations from layer 11; reports token-level FVU, mean L_0 , and dead-feature count.
M3 (Downstream preservation)	Uses cached ImageNet-1K validation activations with an internal stratified 80/20 split, seed = 42; image-level features are obtained by max-pooling patch features; multinomial logistic regression with <code>lbfgs</code> , $C = 1.0$, <code>max_iter</code> =1000, and <code>class_weight</code> =None; the same split is used for raw and reconstructed activations; preservation is reconstructed accuracy divided by raw-activation accuracy.
M4 (Sparse probing)	Uses cached ImageNet-1K validation activations with an internal stratified 80/20 split, seed = 42; SAE codes are max-pooled to image-level features; features are ranked once by descending ANOVA F-statistic using <code>f_classif</code> ; feature budgets $k \in \{1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024\}$; probes use logistic regression with <code>lbfgs</code> , $C = 1.0$, <code>max_iter</code> =500, and <code>class_weight</code> =None; $k = 128$ is the main reported sparse-probing score.
M5 (Monosemanticity)	Top- N activating validation images per feature with $N = 16$; up to 2048 live features are scored; live features are identified by positive maximum activation; embeddings are computed with an independent cross-backbone embedding model; score uses activation-weighted mean pairwise cosine similarity of L2-normalized embeddings; baseline uses 1000 random pseudo-features, each sampling 16 validation images; normalized score is $(s_{\text{obs}} - s_{\text{baseline}})/(1 - s_{\text{baseline}})$.
M6 (Cross-domain)	Evaluated on iNaturalist, CUB-200, and DTD. For each dataset, activations are extracted from the same frozen backbone at layer 11 and normalized with ImageNet statistics; raw activations, SAE reconstructions, and SAE codes are max-pooled to image-level features. Matched linear probes are trained separately for each representation. Preservation is reconstructed-probe accuracy divided by raw-feature accuracy; sparse SAE-code probing uses $k = 128$.
M7 (Probe-based absorption)	Uses max-pooled SAE codes and ImageNet hierarchy sibling groups with minimum group size = 3; features are ranked by <code>f_classif</code> for class-vs-sibling probes; absorption is detected when $F1(k = 1) < \alpha$ and $F1(k = 4) > \beta$; default thresholds are $\alpha = 0.8$ and $\beta = 0.5$; probes use logistic regression with <code>lbfgs</code> , $C = 1.0$, <code>max_iter</code> =500, and <code>class_weight</code> =None; sensitivity analysis sweeps $\alpha \in \{0.7, 0.8, 0.9\}$ and $\beta \in \{0.3, 0.5, 0.7\}$.

852 B Backbone and Hyperparameter Variance Decomposition

Table 3: Variance decomposition η^2 values for each ViSAEBench metric.

Metric	η_{backbone}^2	η_{hp}^2	η_{resid}^2
M1 Localization	0.996	0.002	0.002
M7 Absorption	0.931	0.023	0.046
M5 Monosemanticity	0.916	0.050	0.034
M4 Sparse Probing	0.668	0.070	0.262
M3 Preservation	0.622	0.243	0.135
M2 FVU	0.392	0.523	0.086
M6 iNat Preservation	0.311	0.549	0.141

853 C Probe training and normalization details

854 All probe-training choices are matched across backbones to avoid introducing evaluation bias. For
855 SAE training, all backbones use the same BatchTopK architecture and the same sweep protocol:
856 expansion factor in $\{8, 16, 32\}$, sparsity $K \in \{64, 128, 192, 256\}$, learning rate 10^{-3} , batch size
857 8192, one training epoch, and seed 42. The primary regularizer is the explicit BatchTopK sparsity
858 constraint, so the sparsity regularization strength is harmonized by construction across backbones.

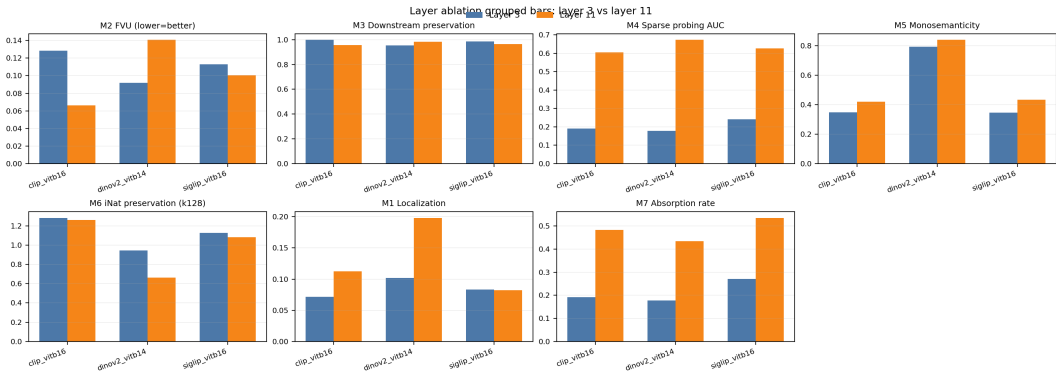


Figure 4: **Layer ablation comparing early-layer and late-layer SAEs.** We compare SAEs trained on layer 3 and layer 11 across seven ViSAEBench metrics for CLIP, DINOv2, and SigLIP. The comparison shows that some metrics are relatively stable across layers, while others change substantially with layer depth. Reconstruction error (M2) and downstream preservation (M3) do not exhibit a single uniform direction across backbones, suggesting that the effect of layer choice is backbone dependent. In contrast, sparse probing performance (M4) and feature absorption (M7) are consistently higher for layer 11, indicating that later-layer SAE features better support sparse concept readout and show stronger absorption structure. Monosemanticity (M5) changes only mildly, preserving the same broad backbone ordering, while localization (M1) increases most clearly for DINOv2 at layer 11.

859 For linear probing in M3 and related probe-based metrics, we use logistic regression with the L-BFGS
 860 solver and $C = 1.0$. In scikit-learn this corresponds to L2-regularized logistic regression under the
 861 default penalty. The same solver, regularization strength, and optimization settings are used for all
 862 backbones and tasks. We do not use class weighting in any probe; equivalently, `class_weight=None`,
 863 the scikit-learn default. Thus, differences in preservation or probing accuracy are not attributable to
 864 backbone-specific probe regularization or class reweighting.

865 Feature normalization is also applied consistently. Before SAE encoding and evaluation, activations
 866 are standardized using the cached dataset statistics: we subtract the activation mean and divide by
 867 the stored standard deviation. The same normalization procedure is used for every backbone, SAE
 868 configuration, and evaluation metric.

869 We did not evaluate whitening-based preprocessing, such as PCA or ZCA whitening, before SAE
 870 training or before probing. Whitening may change the geometry of both the SAE training objective
 871 and the linear probe problem, and we leave this as a follow-up ablation rather than including untested
 872 whitening claims in the present study.

873 D Layer ablation

874 Our main experiments train SAEs on the canonical late-layer representation of each backbone. To test
 875 whether the reported SAE properties are specific to this layer choice, we additionally train matched
 876 SAEs on an earlier transformer layer and compare them against the layer used in the main benchmark.
 877 Specifically, we compare SAEs trained on layer 3 with SAEs trained on layer 11 for CLIP, DINOv2,
 878 and SigLIP, using the same evaluation protocol as in the main paper. This ablation asks whether
 879 properties such as reconstruction fidelity, sparsity, downstream preservation, spatial localization,
 880 monosemanticity, cross-domain preservation, and absorption change systematically as a function of
 881 layer depth.

882 Figure 4 shows that layer choice affects SAE behavior, but not all metrics respond in the same way.
 883 For reconstruction fidelity, measured by M2 FVU, the early and late layers do not follow a universal
 884 trend: CLIP and SigLIP obtain lower FVU at layer 11, whereas DINOv2 obtains lower FVU at
 885 layer 3. This suggests that reconstruction difficulty is not determined by depth alone, but also by the
 886 geometry of the underlying backbone representation at that layer.

887 Downstream preservation, measured by M3, is comparatively stable across the two layers. All three
 888 backbones remain close to the raw-feature reference, and the layer-3 versus layer-11 differences are

889 small relative to the larger cross-backbone effects observed in the main benchmark. This supports the
890 interpretation that SAE reconstruction can preserve much of the downstream signal at both early and
891 late layers, even when other interpretability metrics differ.

892 The clearest systematic layer effect appears in sparse probing and absorption. M4 sparse probing
893 AUC is substantially higher for layer 11 than layer 3 across all three backbones, indicating that
894 late-layer SAE features are more directly useful for sparse concept readout. M7 absorption rate also
895 increases consistently at layer 11, suggesting that later representations exhibit stronger feature-sharing
896 or absorption structure under our criterion. These two trends are consistent with the intuition that
897 later transformer layers contain more semantically organized features than earlier layers.

898 For M5 monosemanticity, the absolute values change only modestly between layers, and the backbone
899 ordering is largely preserved: DINOv2 remains substantially higher than CLIP and SigLIP. Thus,
900 although layer depth affects sparse probing and absorption, the relative monosemanticity differences
901 across backbones are not explained away by the choice of layer. Finally, M1 localization shows a
902 notable increase for DINOv2 at layer 11, while CLIP increases more moderately and SigLIP is nearly
903 unchanged. This again suggests that the spatial organization of SAE features depends jointly on the
904 backbone and the layer at which the SAE is trained.

905 Overall, this ablation supports using a fixed canonical late layer for the main controlled study, while
906 also showing that layer choice is a meaningful experimental axis. The main conclusions are not
907 simply artifacts of evaluating a single layer: backbone-dependent differences persist across layers,
908 but late-layer SAEs tend to produce stronger sparse probing performance and higher absorption rates.

909 E BatchTopK inference threshold

910 BatchTopK training applies sparsity by retaining the top $B \cdot K$ activations across a batch of B
911 tokens. At inference time, however, evaluation should not depend on the other samples that happen to
912 appear in the same batch. We therefore replace batch-coupled top- k selection with a fixed activation
913 threshold estimated after training.

914 After each SAE is trained, we compute a robust inference threshold on held-out validation activations
915 and serialize it with the run artifacts. Specifically, the threshold is saved both in the model check-
916 point, under `_running_threshold`, and in the corresponding `config.yaml`, under `threshold`.
917 This explicit serialization is necessary because the underlying Overcomplete implementation stores
918 `running_threshold` as a Python attribute rather than as a registered PyTorch buffer, so it would
919 not otherwise be captured by `state_dict`.

920 At evaluation time, all metric scripts restore the saved threshold before encoding activations. Thus,
921 SAE inference is deterministic and reproducible from the released checkpoint and configuration alone.
922 We verified that all 60 checkpoints in the main 1M-activation sweep include `_running_threshold`
923 in `sae.pt` and that all corresponding run configurations include the explicit `threshold` field.

924 F M1 Null-Check Statistic

925 For completeness, the SAE-level null-check statistic used for M1 is

$$z = \frac{\bar{I}_{\text{SAE}} - \mathbb{E}[I]}{\widehat{\text{SE}}}, \quad \widehat{\text{SE}} = \frac{\text{sd}(\{\bar{I}_i\}_{i \in \mathcal{E}})}{\sqrt{N_{\text{eval}}}},$$

926 where \mathcal{E} is the set of evaluable features and $N_{\text{eval}} = |\mathcal{E}|$. Because N_{eval} is on the order of 10^4 , even
927 small deviations from the analytic null can produce large z -values; throughout the paper we therefore
928 use z only to test null consistency and use mean Moran’s I for cross-backbone comparison.

929 G Raw-Backbone Spatial Coherence Diagnostic

930 Finding 1 shows that MAE-derived SAE features are near the spatial-random null under M1, whereas
931 the other four backbones produce strongly spatially structured SAE features. To test whether this
932 gap is introduced by SAE training or inherited from the pretrained representation, we repeat the
933 Moran’s- I analysis directly on raw pre-SAE backbone activations.

934 For each backbone, we use the cached ImageNet-1K validation activations from the same layer used
 935 for SAE training. Each raw activation channel is treated as a spatial feature map over the ViT patch
 936 grid. We compute Moran’s I using the same patch-grid statistic as M1, with 8-neighbor connectivity
 937 and the same analytic spatial-random null expectation,

$$\mathbb{E}[I_{\text{null}}] = -\frac{1}{P-1},$$

938 where P is the number of patches. Thus, the null expectation is -0.0051 for patch-16 backbones
 939 with $P = 196$, and -0.0039 for DINOv2-B/14 with $P = 256$.

Table 4: Raw-backbone spatial coherence diagnostic. Moran’s I is computed directly on pre-SAE backbone activation channels using the same patch-grid statistic as M1.

Backbone	Mean raw I	Null $\mathbb{E}[I]$	Raw z	Conclusion
DINOv2	0.3829	-0.0039	202.03	strongly non-random
CLIP	0.2065	-0.0051	196.51	strongly non-random
SigLIP	0.1616	-0.0051	219.64	strongly non-random
DeiT	0.2225	-0.0051	168.08	strongly non-random
MAE	-0.0052	-0.0051	-6.17	near-null effect size

940 The raw-backbone diagnostic mirrors the SAE-level pattern. DINOv2, CLIP, SigLIP, and DeiT all
 941 have large positive raw Moran’s- I values and reject the spatial-random null by a large margin. MAE,
 942 in contrast, has mean raw Moran’s $I = -0.0052$, numerically almost identical to the analytic null
 943 expectation of -0.0051 . Although its raw z -score is below the $|z| < 2$ null-indistinguishability
 944 threshold, the effect size is essentially at the null. This supports the interpretation that MAE’s SAE-
 945 level spatial flatness is inherited from the pretrained backbone representation rather than introduced
 946 by the BatchTopK SAE.

947 **Sensitivity of M1 to design choices.** M1 has three design choices that warrant brief comment.
 948 First, the evaluability threshold of 5 active patches per image was chosen to ensure sufficient spatial
 949 support for Moran’s I to be meaningful while admitting features that activate on small object regions;
 950 the cross-backbone ranking is stable to thresholds in $\{1, 5, 20\}$ on the panel we tested. Second, we
 951 use 8-connected (queen contiguity) neighbours, the standard choice in spatial statistics; the (N/S)
 952 normalization in Moran’s I absorbs the row-sum difference between rook and queen contiguity, and
 953 on blob-shaped ViT features the two yield rank-correlated SAE scores. Third, DINOv2 ViT-B/14
 954 produces a 16×16 patch grid ($P=256$) while the other four backbones produce 14×14 ($P=196$).
 955 The null expectation $\mathbb{E}[I] = -1/(P-1)$ differs by 0.0012 between these two grids, two orders of
 956 magnitude smaller than the cross-backbone effect sizes we observe (DINOv2 leads the next-best
 957 $P=196$ backbone, DeiT, by 0.083 in mean Moran’s I). The raw-backbone diagnostic in this appendix
 958 shows the same DINOv2-leading ranking on raw pre-SAE activations, indicating that the spatial-
 959 coherence gap is not an artifact of the M1 statistic interacting with SAE encoding at different grid
 960 resolutions. A resampling-based control that interpolates DINOv2 feature maps to 14×14 before
 961 computing I would further isolate the effect of pretraining from grid resolution; we leave this to
 962 future work.

963 **H Qualitative audit of M1 spatial coherence**

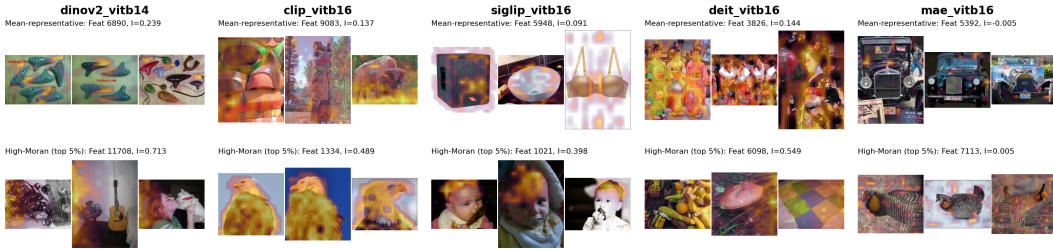


Figure 5: **Qualitative audit of M1 spatial coherence across backbones.** Each column corresponds to one pretrained backbone and shows two SAE features from the BatchTopK $16 \times$, $K = 192$ configuration. The first row shows a mean-representative feature, selected as the evaluable feature whose per-feature Moran’s I is closest to the backbone mean. The second row shows a high-Moran feature from the top 5% of the same distribution. Each cell displays top-activating validation images with patch-level activation overlays. DINOv2, CLIP, SigLIP, and DeiT show spatially contiguous activations that often align with localized visual structure, while MAE remains close to the spatial-random null even in the high-Moran row. This figure illustrates that M1 captures spatial organization, not semantic purity.

964 M1 is designed to measure whether an SAE feature activates in a spatially contiguous region of an
 965 image. It is not, by itself, a semantic-purity metric: a feature can have high spatial autocorrelation
 966 because it activates on a compact object part, a localized texture, a background region, or another
 967 spatially coherent but semantically broad pattern. We therefore include a qualitative audit to show
 968 what different regions of the per-feature Moran’s-I distribution look like across backbones.

969 Figure 5 shows representative SAE features from the BatchTopK $16 \times$, $K = 192$ configuration for
 970 each backbone. For each backbone, we show two features. The first row selects a *mean-representative*
 971 feature, defined as the evaluable feature whose per-feature Moran’s I is closest to the backbone’s mean
 972 Moran’s I for that configuration. This row visualizes what a typical spatial-coherence feature looks
 973 like for the backbone. The second row selects a *high-Moran* feature from the top 5% of the same
 974 per-feature Moran’s-I distribution, visualizing the upper tail of spatially coherent features. For each
 975 selected feature, we display its top-activating validation images with the corresponding patch-level
 976 activation overlay.

977 The qualitative pattern matches the quantitative M1 separation in Section 5.1. DINOv2, CLIP, SigLIP,
 978 and DeiT produce features whose activations often concentrate on contiguous image regions. In
 979 several cases these regions align with visually interpretable content, such as object parts, foreground
 980 objects, or localized textures. This supports the use of M1 as a filter for finding spatially organized,
 981 concept-like SAE features. However, the examples also show why M1 should not be interpreted as a
 982 standalone monosemanticity score: spatially coherent features may still correspond to broad visual
 983 structure rather than a single class-pure concept.

984 MAE provides the clearest contrast. Both the mean-representative MAE feature and the top-5%
 985 MAE feature remain near the spatial-random null, with per-feature Moran’s I close to -0.005 .
 986 Their activation overlays are diffuse and fragmented compared with the other backbones. This
 987 qualitative behavior is consistent with the aggregate result that MAE-derived SAEs are statistically
 988 indistinguishable from the spatial-random null under M1, whereas the other four backbones produce
 989 strongly spatially structured SAE features.

990 Overall, the grid illustrates the intended interpretation of M1: high Moran’s I identifies features with
 991 spatially organized activation patterns and can help surface candidate localized concepts, but semantic
 992 meaningfulness must be assessed jointly with concept-detection metrics such as M4/M5 and with
 993 qualitative inspection.

994 I M4 Sparse Probing Curves

995 In the main text, M4 is summarized by sparse-probing AUC at a fixed feature budget $k = 128$.
996 Here we report the full sparse-probing curves across SAE training sparsity levels and expansion
997 factors. For each SAE, image-level SAE codes are obtained by max-pooling patch-level feature
998 activations. For each ImageNet class, features are ranked by a univariate class-discrimination score
999 on the probe-training split, and logistic-regression probes are trained using the top-ranked features.
1000 We report mean AUC across classes.

1001 Appendix I shows the resulting curves for all five backbones. The cross-backbone ordering is broadly
1002 stable across expansion factors and SAE sparsity levels: DINOv2 and SigLIP consistently provide
1003 the strongest sparse concept readout, CLIP is intermediate, and MAE and DeiT are weaker. This
1004 supports the use of $k = 128$ as a representative summary point in the main text, while also showing
1005 that the M4 trends are not an artifact of a single SAE configuration.

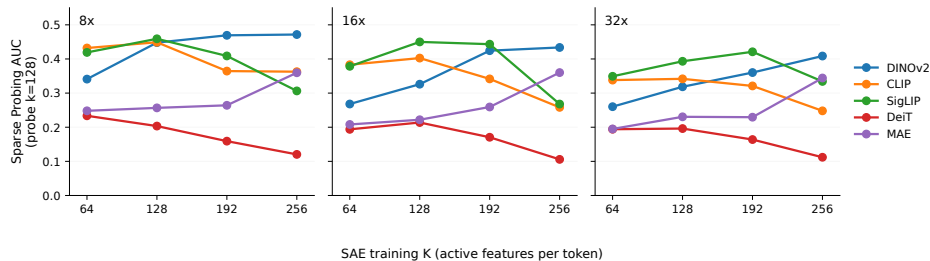


Figure 6: **Sparse probing curves.** ImageNet sparse-probing AUC across SAE training sparsity K and expansion factors. DINOv2 and SigLIP consistently support stronger sparse concept readout than MAE and DeiT across the sweep.

1006 J Per-Task Metric Correlation Panels

1007 The main text summarizes metric–task relationships with a heatmap. Here we provide the same
1008 cross-backbone Spearman correlations as per-task panels, which make it easier to compare which
1009 ViSAEBench metrics are most predictive for each downstream task.

1010 Figure 7 reports correlations between the seven ViSAEBench metrics and four downstream tasks:
1011 retrieval, CUB-200 fine-grained classification, DTD texture classification, and ADE20K segmentation.
1012 M4 sparse probing is the strongest predictor for the two natural-image concept-recognition tasks:
1013 it reaches $\rho = 0.80$ on CUB-200 and $\rho = 1.00$ on DTD. This supports the interpretation that
1014 ImageNet sparse probing captures concept-discriminative structure that transfers to related natural-
1015 image domains. In contrast, M1 Moran’s I is most predictive for segmentation, reaching $\rho = 0.90$,
1016 consistent with the role of spatial coherence in dense prediction. M2 FVU is negatively correlated with
1017 most downstream tasks, reinforcing that reconstruction error alone is not a reliable cross-backbone
1018 quality criterion.

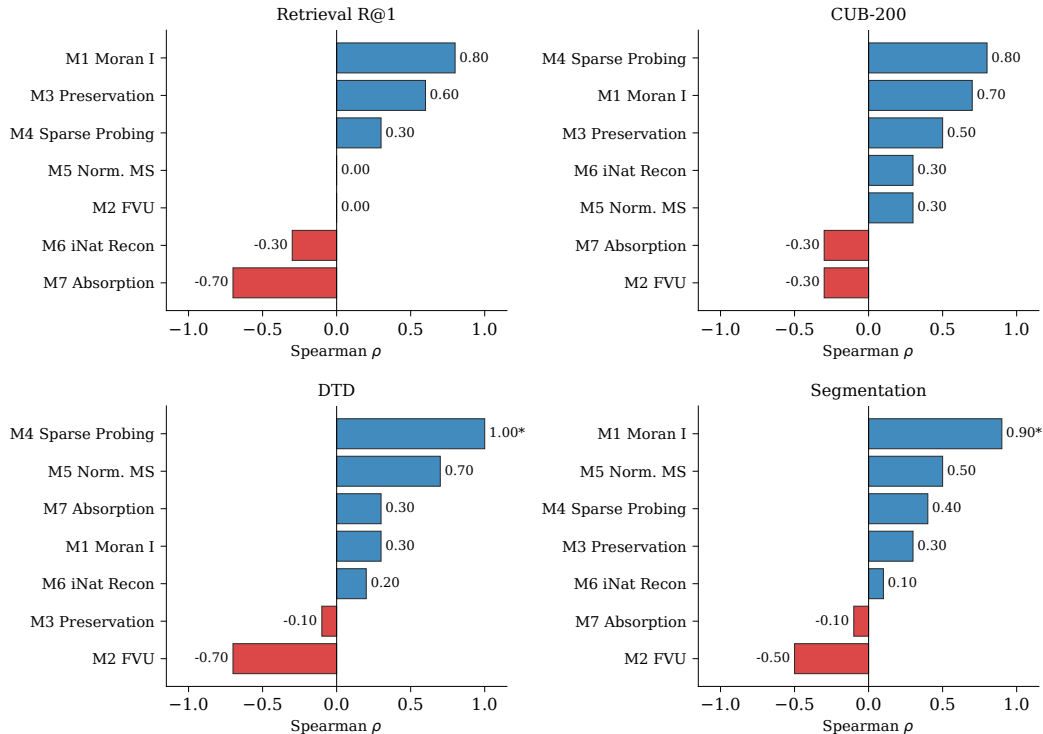


Figure 7: **Per-task metric correlations.** Cross-backbone Spearman correlations between ViSAEBench metrics and downstream tasks. M4 sparse probing best predicts CUB-200 and DTD performance, while M1 spatial coherence best predicts segmentation.

1019 K M5 monosemanticity protocol

Table 5: Exact cross-backbone pairing used for M5 monosemanticity evaluation. Each SAE is evaluated using image embeddings from a backbone different from the one on which the SAE was trained, avoiding circularity in the top-activating-image similarity score.

SAE training backbone	M5 evaluation embedding backbone
DINOv2-B/14	CLIP-ViT-B/16
CLIP-ViT-B/16	DINOv2-B/14
SigLIP-ViT-B/16	DINOv2-B/14
DeiT-ViT-B/16	CLIP-ViT-B/16
MAE-ViT-B/16	CLIP-ViT-B/16

1020 **M5 implementation and normalization.** For each SAE, we encode all ImageNet validation
 1021 activations and obtain one image-level SAE code vector per image by max-pooling over patches, i.e.,
 1022 for each feature we keep its maximum patch activation within the image. A feature is considered
 1023 *live* if its maximum image-level activation over the validation set is strictly positive; features with
 1024 maximum activation equal to zero are treated as dead and excluded from scoring. From the live
 1025 features, we evaluate up to 2,048 features, sampled uniformly without replacement when more than
 1026 2,048 are live and otherwise using all live features.

1027 For each selected feature j , we rank validation images by that feature’s image-level activation and
 1028 take the top- k images, with $k = 16$ unless the dataset slice contains fewer than 16 images. Let $a_{i,j}$
 1029 denote the activation of feature j on top image i . We convert activations to nonnegative weights by

1030 min-max normalizing feature j over all validation images:

$$w_{ij} = \frac{a_{ij} - \min_n a_{nj}}{\max_n a_{nj} - \min_n a_{nj}}.$$

1031 The selected top images are then embedded with the designated cross-backbone evaluator from
 1032 Table 5, and all embeddings are L2-normalized. The feature-level monosemanticity score is the
 1033 activation-weighted mean pairwise cosine similarity

$$\text{MS}_j = \frac{\sum_{u < v} w_{uj} w_{vj} \cos(e_u, e_v)}{\sum_{u < v} w_{uj} w_{vj}},$$

1034 computed only when the denominator is nonzero. The raw M5 score, M5_{raw} , is the mean of MS_j
 1035 over all scored features.

1036 Raw M5 reflects two factors: (i) SAE-dependent image selection, i.e., which images a feature
 1037 activates on, and (ii) evaluator-dependent embedding geometry, i.e., how tightly those selected images
 1038 cluster in the evaluator’s representation space. Even if image sets were sampled at random, different
 1039 evaluation backbones can induce different expected pairwise cosine similarities because of anisotropy,
 1040 class structure, and feature scaling in their embedding spaces. Raw M5 is therefore not a pure SAE
 1041 property.

1042 To correct for this evaluator-specific similarity floor, we estimate a backbone-specific geometric
 1043 baseline using random pseudo-features. Concretely, we sample 1,000 random image groups, matching
 1044 the same top- k size used for real features, from the exact embedded image pool used in scoring.
 1045 Each pseudo-feature uses uniform pair weights, equivalently the unweighted mean pairwise cosine
 1046 similarity. The baseline $\text{M5}_{\text{baseline}}$ is the average over these 1,000 pseudo-features and measures the
 1047 expected similarity induced by the evaluator’s embedding geometry alone.

1048 We report normalized M5 as

$$\text{M5}_{\text{norm}} = \frac{\text{M5}_{\text{raw}} - \text{M5}_{\text{baseline}}}{1 - \text{M5}_{\text{baseline}}}.$$

1049 This rescales the baseline-corrected score by the available cosine-similarity headroom above baseline,
 1050 so normalized values can be interpreted as the fraction of possible improvement over the evaluator-
 1051 specific random baseline. As a result, cross-backbone rankings can differ between raw and normalized
 1052 M5: raw rankings include both SAE signal and evaluator-geometry bias, whereas normalized rankings
 1053 more directly reflect how far SAE-selected image sets rise above evaluator-specific chance structure.

1054 **Within-class vs. between-class decomposition of M5.** To test whether high M5 reflects strict class-
 1055 pure selectivity or broader visual clustering, we decompose each feature’s pairwise similarity into
 1056 within-class and between-class terms using ImageNet validation labels. For a scored feature j with
 1057 top- k images, let embeddings be $\{e_i\}_{i=1}^k$, labels $\{y_i\}_{i=1}^k$, and activation-derived weights $\{w_{ij}\}_{i=1}^k$
 1058 as defined above. Define pair weight $\omega_{uv} = w_{uj} w_{vj}$ and cosine similarity $c_{uv} = \cos(e_u, e_v)$ for $u < v$.

1059 We partition pairs into

$$\mathcal{W}_j = \{(u, v) : y_u = y_v, u < v\}, \quad \mathcal{B}_j = \{(u, v) : y_u \neq y_v, u < v\},$$

1060 and compute

$$\text{MS}_j^{\text{within}} = \frac{\sum_{(u,v) \in \mathcal{W}_j} \omega_{uv} c_{uv}}{\sum_{(u,v) \in \mathcal{W}_j} \omega_{uv}}, \quad \text{MS}_j^{\text{between}} = \frac{\sum_{(u,v) \in \mathcal{B}_j} \omega_{uv} c_{uv}}{\sum_{(u,v) \in \mathcal{B}_j} \omega_{uv}},$$

1061 whenever the corresponding denominator is nonzero. We also compute the pair-mass fractions

$$\rho_j^{\text{within}} = \frac{\sum_{(u,v) \in \mathcal{W}_j} \omega_{uv}}{\sum_{u < v} \omega_{uv}}, \quad \rho_j^{\text{between}} = 1 - \rho_j^{\text{within}},$$

1062 which yield the exact convex decomposition

$$\text{MS}_j = \rho_j^{\text{within}} \text{MS}_j^{\text{within}} + \rho_j^{\text{between}} \text{MS}_j^{\text{between}}.$$

1063 At the aggregate level, averaging over scored features, high-M5 features typically have elevated
 1064 $\text{MS}^{\text{within}}$ and nontrivial $\text{MS}^{\text{between}}$. Thus, their top-activating images cluster strongly within class
 1065 but can also remain highly similar across different classes. This supports the claim in §5.6 that M5
 1066 captures within-feature embedding-space consistency rather than class-pure selectivity or fine-grained
 1067 discriminability.

1068 **L M6 cross-domain protocol**

1069 For each out-of-domain dataset, we train probes using only that dataset’s own training split and
 1070 evaluate on its held-out test split. No ImageNet-trained classifier is reused for M6. Instead, each
 1071 representation is probed independently: raw backbone activations, SAE-reconstructed activations,
 1072 and sparse SAE feature vectors each receive their own classifier trained from scratch under the
 1073 same hyperparameters. For the CPU implementation we use scikit-learn logistic regression with
 1074 L-BFGS, $C = 1.0$, and `max_iter = 500`; for the GPU implementation we use Adam with learning
 1075 rate 10^{-2} , weight decay 10^{-4} , 50 epochs, and batch size 2048. The M6 preservation ratio is then
 1076 computed as reconstructed-probe accuracy divided by raw-activation-probe accuracy on the same
 1077 OOD dataset. This protocol makes the ratio an estimate of SAE-induced information loss under
 1078 matched OOD supervision, rather than a measure of how well an ImageNet classifier transfers to the
 1079 OOD label space. In practice the ratio is usually at or below 1, with values near 1 indicating that SAE
 1080 reconstruction preserves nearly all task-relevant OOD information available in the raw representation.
 1081 This protocol isolates SAE-induced loss rather than backbone capacity or ImageNet-to-OOD classifier
 1082 transfer.

1083 **M Full Cross-Domain Probe Accuracies (M6)**

1084 This appendix reports the raw, reconstructed, and preserved linear-probe accuracies underlying the M6
 1085 (cross-domain preservation) score for all five backbones across the BatchTopK SAE sweep ($n=60$:
 1086 $\{8\times, 16\times, 32\times\} \times \{k=64, 128, 192, 256\}$). For each dataset $\mathcal{D} \in \{\text{iNaturalist, CUB-200, DTD}\}$ we
 1087 train matched linear probes on (i) raw backbone features $\text{Acc}_{\text{raw}}^{\mathcal{D}}$ and (ii) SAE-reconstructed features
 1088 $\text{Acc}_{\text{recon}}^{\mathcal{D}}$, and report the preservation ratio $P^{\mathcal{D}} = \text{Acc}_{\text{recon}}^{\mathcal{D}} / \text{Acc}_{\text{raw}}^{\mathcal{D}}$.

1089 Table 6 aggregates over the 12 SAE configurations per backbone (reporting mean \pm std across the
 1090 sweep). $\text{Acc}_{\text{raw}}^{\mathcal{D}}$ is a property of the backbone alone and is therefore reported once per row.

Table 6: Cross-domain probe accuracies and preservation ratios. Recon and P columns report mean \pm std across $n=12$ BatchTopK SAE configurations per backbone; raw is the backbone-only probe.

Backbone	iNaturalist (10K classes)			CUB-200 (200 classes)			DTD (47 classes)		
	Raw	Recon	P	Raw	Recon	P	Raw	Recon	P
CLIP-B/16	0.093	0.074 \pm 0.007	0.793 \pm 0.078	0.575	0.517 \pm 0.018	0.900 \pm 0.032	0.793	0.773 \pm 0.014	0.974 \pm 0.018
SigLIP-B/16	0.119	0.101 \pm 0.006	0.852 \pm 0.050	0.693	0.658 \pm 0.011	0.949 \pm 0.016	0.836	0.825 \pm 0.006	0.987 \pm 0.007
DINOv2-B/14	0.235	0.210 \pm 0.009	0.893 \pm 0.037	0.831	0.817 \pm 0.007	0.984 \pm 0.008	0.831	0.816 \pm 0.005	0.983 \pm 0.007
DeiT-B/16	0.074	0.057 \pm 0.008	0.770 \pm 0.107	0.623	0.528 \pm 0.052	0.847 \pm 0.084	0.699	0.672 \pm 0.009	0.961 \pm 0.012
MAE-B/16	0.037	0.032 \pm 0.002	0.864 \pm 0.062	0.420	0.381 \pm 0.020	0.908 \pm 0.048	0.691	0.661 \pm 0.014	0.957 \pm 0.021

1091 **Observations.** (i) Raw probe accuracy varies by an order of magnitude across datasets and back-
 1092 bones (e.g. MAE-B/16 reaches 3.7% on iNaturalist-10K vs. 69.1% on DTD), which is exactly why
 1093 $P^{\mathcal{D}}$ rather than raw $\text{Acc}_{\text{recon}}$ is the correct preservation signal. (ii) Preservation is consistently highest
 1094 on DTD ($P \geq 0.957$ for all backbones) and lowest on iNaturalist, where the long-tailed 10K-way
 1095 label space amplifies any reconstruction error. (iii) DeiT-B/16 shows the largest variance across
 1096 SAE configurations (e.g. $P^{\text{CUB}} = 0.847 \pm 0.084$), indicating its representations are more sensitive
 1097 to the choice of SAE width and sparsity than the other backbones; DINOv2-B/14 is the most stable
 1098 ($P^{\text{CUB}} = 0.984 \pm 0.008$).

1099 **N M7 threshold sensitivity**

1100 We evaluate the sensitivity of M7 absorption to the overlap threshold α and size-ratio threshold β by
 1101 sweeping $\alpha \in \{0.7, 0.8, 0.9\}$ and $\beta \in \{0.3, 0.5, 0.7\}$. Table 7 reports the per-backbone absorption
 1102 rates. DeiT has the lowest absorption in 7 of 9 settings. The two exceptions occur at the strictest
 1103 overlap threshold, $\alpha = 0.9$, with $\beta \in \{0.5, 0.7\}$, where MAE becomes lower than DeiT. This
 1104 confirms that the DeiT absorption paradox is not an artifact of the default threshold choice, while
 1105 also showing that the exact rank ordering is threshold-sensitive in the high-stringency regime.

Table 7: Sensitivity of M7 absorption rate to the absorption thresholds α and β . Lower is conventionally better. Bold indicates the lowest-absorption backbone for each threshold setting.

α	β	DINOv2	CLIP	SigLIP	DeiT	MAE
0.7	0.3	0.385	0.562	0.502	0.155	0.521
0.7	0.5	0.294	0.404	0.408	0.121	0.291
0.7	0.7	0.098	0.200	0.208	0.038	0.106
0.8	0.3	0.513	0.653	0.642	0.332	0.645
0.8	0.5	0.423	0.494	0.547	0.298	0.415
0.8	0.7	0.223	0.291	0.347	0.211	0.226
0.9	0.3	0.725	0.796	0.796	0.604	0.785
0.9	0.5	0.634	0.638	0.702	0.570	0.555
0.9	0.7	0.434	0.434	0.502	0.483	0.366

1106 O The Codes-lift Inversion

1107 Across the four downstream tasks where we evaluate both raw features and SAE codes (retrieval R@1,
 1108 CUB-200, DTD, segmentation), we observe a consistent and counterintuitive pattern: backbones with
 1109 weaker raw-feature performance gain proportionally more from SAE decomposition than backbones
 1110 with stronger raw-feature performance. We refer to this as the *codes-lift inversion*.

1111 **Setup** For each (backbone, task) pair, we compute the lift ratio $L =$
 1112 $\text{perf}(\text{SAE codes})/\text{perf}(\text{raw features})$, averaged across the 12 SAE configurations per back-
 1113 bone. We then compute the Spearman rank correlation between mean raw performance and mean lift
 1114 across the five backbones ($n = 5$).

1115 **Results** On three of four tasks, the correlation is strongly negative (Table 8).

Task	$\rho(\text{raw}, L)$	p	Range of lifts
Retrieval R@1	-0.900	0.037	DeiT 0.99 to MAE 1.18
DTD	-0.900	0.037	SigLIP 1.02 to DeiT 1.06
CUB-200	-0.700	0.188	DINOv2 0.99 to MAE 1.19
Segmentation mIoU	+0.500	0.391	DeiT 0.43 to DINOv2 0.75

Table 8: Cross-backbone Spearman correlation between raw-feature performance and SAE codes lift. Negative correlations indicate the inversion: weaker raw backbones gain more from SAE decomposition. Segmentation is the lone exception.

1116 The magnitude of the inversion is substantial. MAE, the worst raw performer on retrieval (R@1
 1117 = 0.335), exhibits a 17.7% codes lift, while DeiT, the best raw retrieval performer (R@1 = 0.692),
 1118 shows a slight codes regression (lift = 0.988). On CUB-200, MAE gains 19.4% from SAE de-
 1119 composition while DINOv2 loses 1.4%. The pattern is stable across all 12 configurations within
 1120 each backbone: per-config lift standard deviations are small (between 0.006 for DTD on DeiT and
 1121 0.094 for segmentation on MAE), indicating the effect is a property of the backbone rather than
 1122 configuration noise.

1123 Segmentation is the lone exception. The correlation flips sign to +0.50 and no backbone reaches
 1124 a lift above 1.0 (every backbone loses mIoU under SAE decomposition, with DINOv2 retaining
 1125 the largest fraction at 0.75 and DeiT the smallest at 0.43). This dissociation matches the broader
 1126 spatial-coherence finding from Section 5.1: the spatial structure is a property of the backbone that
 1127 SAE decomposition cannot manufacture, only preserve to varying degrees.

1128 **Interpretation** We interpret the inversion as a ceiling effect on SAE-extractable concept structure
 1129 rather than as evidence that weaker backbones are secretly better. A backbone with raw representations
 1130 that already align well with task labels has less headroom for the SAE’s sparse decomposition to
 1131 add discriminative structure on top. A backbone with weaker raw alignment offers more room for
 1132 the sparse code to contribute task-relevant axes, and the BatchTopK reconstruction loss appears to

1133 capture some of that available signal. The inversion does not hold for spatial tasks because the SAE’s
1134 sparse codes are reconstructed from the backbone’s feature space, and patch-level spatial coherence
1135 is a geometric property of that space which decomposition cannot create.

1136 **Practical implications** The inversion has two consequences for practitioners. First, the relative
1137 gap between backbones narrows substantially after SAE decomposition on retrieval, fine-grained
1138 classification, and texture, which weakens the case for selecting a backbone purely on raw downstream
1139 performance when the downstream consumer is the SAE codes themselves. Second, the inversion does
1140 not generalize to spatial tasks; for segmentation and dense prediction, raw-feature ranking is preserved
1141 through SAE decomposition and DINOv2 retains its dominance regardless of lift considerations.

1142 **Limitations** The result is computed across $n = 5$ backbones, so the ρ values should be read as
1143 suggestive of an effect rather than tight estimates of its magnitude. Per-backbone variance across
1144 the 12 SAE configurations is small (typically below 0.05 in lift units), but the cross-backbone trend
1145 would benefit from replication on additional pretraining paradigms before being treated as a general
1146 law. We do not have a mechanistic account of why the inversion is task-dependent in the specific
1147 pattern observed, and we leave a feature-level intervention study to future work.

1148 P Practitioner guide for backbone selection

1149 This appendix condenses the empirical task-to-backbone mapping into a single-page reference.
1150 Recommendations are derived from per-backbone means across 12 SAE configurations on the 1.28M
1151 ImageNet training split, and the supporting capability metric is the one with the highest cross-
1152 backbone Spearman rank correlation to the task in our pre-registered analysis (Section ??, $n = 5$
1153 backbones).

1154 P.1 Decision summary

1155 Figure 8 provides the top-level routing. In text form:

1156 **Dense spatial prediction.** If the task involves semantic segmentation or dense feature transfer to
1157 spatial tasks, use DINOv2. The supporting metric is M1 (Moran’s I), which correlates with raw
1158 segmentation mIoU at $\rho = 0.90$, $p = 0.037$ across backbones. DINOv2’s mean raw mIoU on
1159 ADE20K is 0.171 versus 0.058 for the next backbone (CLIP), a $2.9\times$ lead that is preserved through
1160 SAE decomposition (codes mIoU 0.127 versus 0.036).

1161 **Fine-grained concept discrimination.** If the task involves CUB-200, iNaturalist transfer, or object
1162 classification with many similar classes, use DINOv2 for raw features and either DINOv2 or SigLIP
1163 for SAE codes. The supporting metric is M4 (Sparse Probing AUC), which correlates with iNat at
1164 $\rho = 0.90$ and CUB-200 at $\rho = 0.80$. DINOv2 leads CUB raw accuracy at 0.831 with SigLIP at
1165 0.693; on SAE codes the gap narrows to 0.819 versus 0.748.

1166 **Texture and material classification.** If the task involves DTD or analogous datasets, use SigLIP.
1167 The supporting metric is again M4, with $\rho = 1.00$ to DTD ($p = 0.017$, $n = 5$). SigLIP’s raw DTD
1168 accuracy is 0.836 with DINOv2 essentially tied at 0.831; on SAE codes SigLIP edges DINOv2 by
1169 0.002. The earlier prediction that masked-pretraining (MAE) would win on texture did not hold.

1170 **Image-image retrieval.** If the task is retrieval (R@1, R@5, R@10) and interpretability is secondary,
1171 use DeiT. The supporting metric is M3 (Preservation), which correlates with retrieval at $\rho = 0.60$.
1172 DeiT achieves R@1 = 0.692 on raw features and 0.684 on SAE codes, well ahead of other backbones.
1173 We note a caveat: DeiT exhibits the worst sparse probing performance (0.418) and a high feature
1174 absorption rate, so applications requiring both retrieval accuracy and feature interpretability should
1175 consider DINOv2 (R@1 = 0.548) or SigLIP (R@1 = 0.526) as more balanced alternatives.

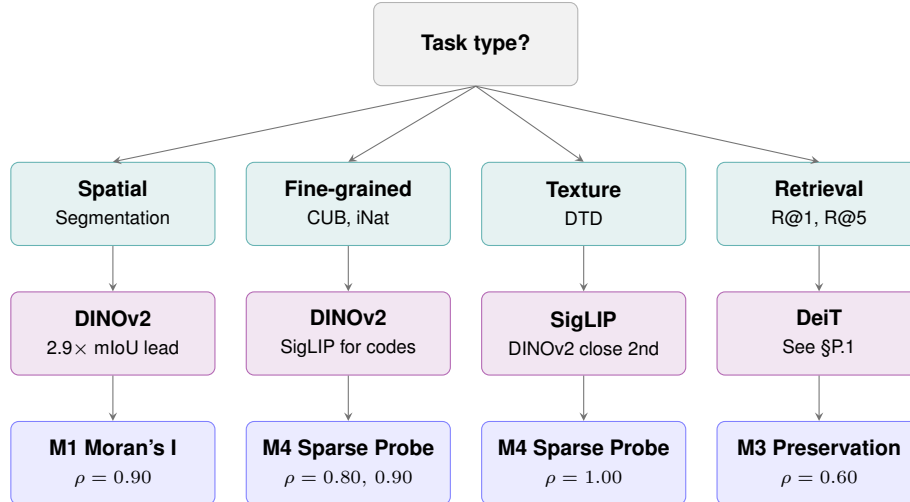


Figure 8: Backbone selection flowchart. Each column maps a task family (teal) to the recommended backbone (purple) and the supporting capability metric (blue), with ρ values from cross-backbone Spearman correlation on $n = 5$ backbones. Caveats and limitations are discussed in the surrounding text.

1176 P.2 Backbones we do not recommend

1177 MAE places last or near-last on every raw downstream task in our suite. We retain it in the benchmark
 1178 as a reference point for the masked-pretraining paradigm and to enable backbone-specific failure-
 1179 mode reporting, but it is not a recommended choice for any of our five task families.

1180 P.3 What not to select on

1181 M2 (Fractional Variance Unexplained) correlates negatively or near-zero with every downstream task
 1182 we evaluate (retrieval 0.00, CUB -0.30 , DTD -0.70 , segmentation -0.50 , iNat -0.60). Selecting
 1183 a backbone on FVU alone is the most consistent way to make a poor choice in our results. M3
 1184 (Preservation) is a usable proxy when M1 or M4 cannot be computed, but it is weaker than either of
 1185 those two metrics on every task except retrieval.

1186 P.4 Codes-lift caveat

1187 Practitioners selecting a backbone for SAE-codes consumption should expect the relative gap between
 1188 backbones to compress on retrieval, CUB, and DTD, but not on segmentation. We document this
 1189 effect in detail in Appendix [Rajamanoharan et al., 2024]. The strongest raw backbone is not always
 1190 the strongest SAE-codes backbone, and the codes ranking should be checked directly against the
 1191 target task before committing to a choice.

1192 P.5 Caveat on sample size

1193 All cross-backbone correlations are computed across $n = 5$ backbones. Three of four pre-registered
 1194 hypotheses pass at $\rho \geq 0.7$ with $p \leq 0.05$, and one (M3 to retrieval) is directionally consistent but
 1195 does not reach significance. The guidance above should be read as evidence-grounded but not as a
 1196 substitute for evaluating SAE features on the practitioner’s actual target task when stakes are high.

1197 Q Architectural robustness: do these findings hold under JumpReLU?

1198 **Setup.** We trained JumpReLU SAEs [Rajamanoharan et al., 2024] on the same five backbones
 1199 at expansion factors $16\times$ and $32\times$ and target L_0 values of 64 and 128, using a uniform recipe
 1200 (Silverman kernel, bandwidth 10^{-3} , adaptive sparsity coefficient with multiplicative controller). Of
 1201 the 20 resulting configurations, 17 settled within ± 2 of their target L_0 and are reported here. Three

1202 configurations (SigLIP $16\times$ at $L_0=128$, SigLIP $32\times$ at $L_0=128$, DINOv2 $32\times$ at $L_0=128$) failed
 1203 to converge to the target L_0 within the one-epoch training budget under the uniform recipe and are
 1204 excluded; we believe these reflect known JumpReLU sensitivity to controller hyperparameters rather
 1205 than fundamental incompatibility, and leave per-backbone tuning to future work.

1206 **Cross-family rank consistency.** Table 9 reports the Spearman rank correlation across the five back-
 1207 bones between BatchTopK and JumpReLU mean values for each of the seven metrics. Correlations
 1208 range from 0.80 (M6) to 1.00 (M2, M5), with six of seven significant at $p < 0.05$ despite the small
 1209 sample size of $n = 5$ backbones. DINOv2 ranks first overall under both SAE families and MAE
 1210 ranks last under both. We interpret this as evidence that the cross-backbone findings reported in the
 1211 main text reflect properties of the underlying vision backbones rather than artifacts of the BatchTopK
 1212 SAE family.

1213 **Absolute differences.** Although ordering is preserved, absolute values differ between SAE fam-
 1214 ilies. JumpReLU exhibits uniformly higher FVU than BatchTopK at matched L_0 (mean increase
 1215 0.037 across backbones), and uniformly lower sparse probing AUC (mean decrease 0.075). These
 1216 differences are consistent across all five backbones, indicating they reflect a property of the SAE
 1217 family rather than an interaction with any specific backbone.

1218 **A note on CLIP.** The single most striking cross-family difference concerns dead-feature rates. At
 1219 matched $16\times$ configurations, CLIP under BatchTopK has 30.6% dead features on average, compared
 1220 to 9.5% under JumpReLU; the other four backbones differ by less than 5 percentage points between
 1221 families. Across the full BatchTopK sweep, CLIP dead rates range from 20.8% to 47.0%, an outlier
 1222 relative to the other backbones (DINOv2: 0.0% to 6.9%; SigLIP: 2.9% to 13.6%). We do not have a
 1223 definitive explanation but note that this interaction between CLIP’s representation and the BatchTopK
 1224 forced-allocation mechanism warrants further investigation.

1225 **Sensitivity check.** Excluding CLIP $32\times$ configurations (where dead-feature rates are highest) from
 1226 the per-backbone aggregation leaves all seven cross-family rank correlations unchanged to three
 1227 decimal places, indicating the robustness conclusion does not depend on the dead-feature outlier.

Table 9: Cross-family rank consistency across the five backbones. Rankings are best to worst per metric (lower-is-better metrics inverted accordingly). Spearman ρ computed on the per-backbone means across all configurations satisfying $|L_0 - \text{target}| \leq 2$. * $p < 0.05$.

Metric	BatchTopK ranking	JumpReLU ranking	ρ	p
M1 Localization	DINOv2, DeiT, CLIP, SigLIP, MAE	DINOv2, CLIP, DeiT, SigLIP, MAE	0.900*	0.037
M2 FVU	CLIP, SigLIP, DINOv2, DeiT, MAE	CLIP, SigLIP, DINOv2, DeiT, MAE	1.000*	0.000
M3 Preservation	DINOv2, DeiT, MAE, SigLIP, CLIP	DINOv2, DeiT, SigLIP, MAE, CLIP	0.900*	0.037
M4 Sparse Probing	SigLIP, DINOv2, CLIP, DeiT, MAE	DINOv2, SigLIP, CLIP, DeiT, MAE	0.900*	0.037
M5 Norm. MS	CLIP, SigLIP, DINOv2, DeiT, MAE	CLIP, SigLIP, DINOv2, DeiT, MAE	1.000*	0.000
M6 iNat	SigLIP, DINOv2, CLIP, MAE, DeiT	DINOv2, SigLIP, CLIP, DeiT, MAE	0.800	0.104
M7 Absorption	DeiT, DINOv2, MAE, SigLIP, CLIP	DeiT, DINOv2, MAE, CLIP, SigLIP	0.900*	0.037