
Same Concept, Different Directions: Cross-Modal Feature Heterogeneity in Sparse Autoencoders

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Vision–language models map images and text into a joint embedding space. How-
2 ever, these embeddings often entangle multiple semantic features, which limits their
3 interpretability and controllability. While sparse autoencoders have emerged as a
4 useful tool for decomposing these embeddings into monosemantic features, their
5 application to joint embedding spaces has largely relied on an implicit, untested as-
6 sumption that semantically corresponding features share the same directions across
7 modalities. In this paper, we challenge this assumption by identifying discrepan-
8 cies in feature directions for the same concept across image and text modalities,
9 a phenomenon we term *cross-modal feature heterogeneity*. We demonstrate that
10 this heterogeneity is a key driver of the *modality split*, where a shared concept
11 activates different latents depending on the modality. This finding further reveals
12 why aligning latent activations alone is insufficient to resolve the underlying feature
13 mismatch. To address this misalignment, we propose an approach that trains sparse
14 autoencoders to preserve the unique feature geometry of each modality and aligns
15 corresponding features post hoc. Our method improves reconstruction fidelity and
16 enhances performance in cross-modal retrieval and concept steering.

17 1 Introduction

18 Vision-language models (VLMs) map images and text into a joint representation space, which supports
19 a wide range of downstream tasks [42, 54, 6] and serves as the foundation for generative VLMs [30,
20 46]. Despite their empirical success, the joint representation remains difficult to interpret [36]. These
21 embeddings are typically *polysemantic*, as they encode multiple semantic concepts in a single vector,
22 and it remains unclear how each concept is encoded and shared across modalities.

23 The linear representation hypothesis [16, 39] provides a useful framework for studying this question,
24 positing that each embedding can be expressed as a linear combination of a small number of
25 *monosemantic* features. Sparse autoencoders (SAEs) operationalize this hypothesis by mapping each
26 embedding to a *sparse latent code*, whose active coordinates indicate which monosemantic features
27 are present in the embedding. Building on their success in unimodal settings [23, 50, 18, 34], recent
28 work has applied SAEs to joint embedding spaces to recover monosemantic features shared across
29 image and text modalities [11, 55, 56, 38, 37, 25, 14].

30 When applied to joint embeddings of VLMs, prior work [25] observes a phenomenon known as
31 *modality split*, where the same concept activates different latent coordinates across modalities. This
32 split makes latent codes difficult to use across modalities, motivating prior approaches to align latent
33 activations across modalities [25, 14]. However, these approaches treat modality split as a mismatch
34 in latent activations, while implicitly assuming that semantically corresponding features share the
35 same directions across modalities.

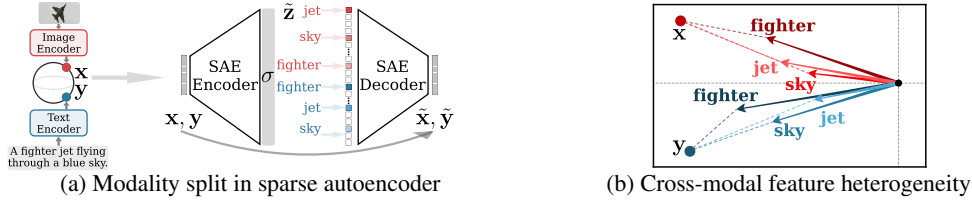


Figure 1: Illustration of modality split and cross-modal feature heterogeneity on joint embedding spaces. (a) An SAE encodes an image or text embedding (\mathbf{x} or \mathbf{y}) from a VLM into a sparse latent code $\hat{\mathbf{z}}$ and then reconstructs the original input as $\hat{\mathbf{x}}$ or $\hat{\mathbf{y}}$. Ideally, the same concept should activate the same coordinate in $\hat{\mathbf{z}}$ across modalities. However, prior work observes the *modality split*, where the same concept (e.g., *sky*, *jet*, or *fighter*) activates different coordinates for images (red) and text (blue). (b) We show that this split is driven by *cross-modal feature heterogeneity*, where corresponding features (shown as arrows) fail to align directionally across modalities within the embedding space.

36 In this paper, we revisit this assumption by asking whether semantically corresponding feature
 37 directions are indeed aligned across modalities in the learned embedding space. We characterize
 38 *cross-modal feature heterogeneity*, a phenomenon in which the same semantic concept can be
 39 represented by different feature directions across image and text modalities, as illustrated in Figure 1.

40 Accounting for this heterogeneity is crucial as it redefines our interpretation of the modality split.
 41 When corresponding image and text features have different directions in the joint embedding space,
 42 an SAE naturally assigns them to different latent coordinates to ensure precise reconstruction.
 43 Consequently, while forcing these latent activations to align across different modalities might mitigate
 44 the modality split, it risks degrading feature recovery by collapsing geometrically distinct directions
 45 into a single coordinate. This motivates our approach, which trains SAEs to preserve the unique
 46 feature geometry of each modality and then aligns the corresponding features post hoc.

47 Our contributions are summarized below.

- 48 • In Section 3, we characterize *cross-modal feature heterogeneity*, where the same concept can have
 49 different feature directions across modalities in the joint embedding space.
- 50 • In Section 4, we show that this heterogeneity can explain *modality split* in multimodal SAEs and
 51 analyze how existing alignment approaches trade reconstruction quality for latent alignment.
- 52 • In Section 5, we propose an approach that preserves the unique feature geometry of each modality
 53 and aligns corresponding latent coordinates post hoc without sacrificing reconstruction quality.
- 54 • In Section 6, we evaluate our approach, showing that preserving reconstruction quality is important
 55 for improving performance on cross-modal tasks, including retrieval and concept steering.

56 Due to space constraints, a detailed comparison with related work is provided in Appendix A.

57 2 Preliminaries

58 We formalize joint embeddings from VLMs as linear combinations of monosemantic features. We
 59 then introduce SAEs as a framework for recovering feature directions from these embeddings.

60 **Embeddings as linear combinations of features.** Following the linear representation hypothesis [3,
 61 16, 39, 12], we assume that VLM embeddings $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ are polysemantic, meaning that each
 62 embedding can be expressed as a linear combination of monosemantic features. A monosemantic
 63 feature is a direction in the joint embedding space that represents a semantically coherent concept.

64 Formally, we assume the existence of a latent code $\mathbf{z} := (z_1 \cdots z_n)^\top \in \mathbb{R}_+^n$ and *feature matrices*
 65 $\Phi := (\phi_1 \cdots \phi_n) \in \mathbb{R}^{d \times n}$ for images and $\Psi := (\psi_1 \cdots \psi_n) \in \mathbb{R}^{d \times n}$ for text. The column
 66 vectors $\phi_i, \psi_i \in \mathbb{R}^d$ denote the *feature direction vectors* of the i -th latent concept for images and
 67 text, respectively, and we assume that they have unit norm.

68 Under the linear representation hypothesis, an image embedding is written as $\mathbf{x} = \sum_{i \in [n]} z_i \phi_i$,
 69 where each activation value $z_i \in \mathbb{R}_+$ indicates the strength of the corresponding feature, whereas
 70 $z_i = 0$ indicates that it is inactive. When an image–text pair is considered, we assume that the two

71 embeddings share the same latent code \mathbf{z} , while their feature directions may differ across modalities:

$$\mathbf{x} = \Phi\mathbf{z} = \sum_{i \in [n]} z_i \phi_i, \quad \mathbf{y} = \Psi\mathbf{z} = \sum_{i \in [n]} z_i \psi_i. \quad (1)$$

72 Unlike prior work [25], we impose no constraint such as $\phi_i = \psi_i$. Thus, we allow $\phi_i \neq \psi_i$,
73 capturing the possibility that the same concept is encoded along different directions across modalities.

74 We further assume that the coordinates of the latent code \mathbf{z} are independent and sparse. Specifically,
75 the sparse factor s is defined as the probability that each coordinate is inactive,

$$\Pr(z_i = 0) = s, \quad i \in [n], \quad (2)$$

76 where s is close to 1. This implies that only a small number of features are active in each embedding.

77 **Sparse autoencoder.** SAEs encode each embedding into a sparse latent code and decode this
78 code to reconstruct the original embedding [23, 50, 18, 44, 55, 12]. Ideally, each active coordinate
79 of the latent code represents the strength of a monosemantic feature in the embedding, while the
80 corresponding decoder column provides its feature direction.

81 Specifically, let $\mathbf{W} := [\mathbf{w}_1 \cdots \mathbf{w}_m] \in \mathbb{R}^{d \times m}$ denote the weight matrix of an SAE. For simplicity, we
82 omit bias terms and use the same weight matrix for the encoder and decoder. The encoder produces
83 an *estimated latent code* $\tilde{\mathbf{z}} := (\tilde{z}_1 \cdots \tilde{z}_m)^\top \in \mathbb{R}_+^m$. For image and text embeddings, we write

$$\tilde{\mathbf{z}}(\mathbf{x}) := \sigma(\mathbf{W}^\top \mathbf{x}), \quad \tilde{\mathbf{z}}(\mathbf{y}) := \sigma(\mathbf{W}^\top \mathbf{y}),$$

84 where σ denotes an activation function¹ that induces sparsity. The decoder then produces

$$\tilde{\mathbf{x}}(\mathbf{x}) := \mathbf{W}\tilde{\mathbf{z}}(\mathbf{x}) = \sum_{j \in [m]} \tilde{z}_j(\mathbf{x}) \mathbf{w}_j, \quad \tilde{\mathbf{y}}(\mathbf{y}) := \mathbf{W}\tilde{\mathbf{z}}(\mathbf{y}) = \sum_{j \in [m]} \tilde{z}_j(\mathbf{y}) \mathbf{w}_j. \quad (3)$$

85 The SAE is trained to reconstruct each embedding using a reconstruction loss. For image embeddings,

$$\mathcal{L}_{\text{rec}}(\mathbf{W}; \Phi) := \mathbb{E}_{\mathbf{x}} \|\mathbf{x} - \tilde{\mathbf{x}}(\mathbf{x})\|_2^2 = \mathbb{E}_{\mathbf{z}} \|\Phi\mathbf{z} - \mathbf{W}\sigma(\mathbf{W}^\top \Phi\mathbf{z})\|_2^2, \quad (4)$$

86 and we use $\mathcal{L}_{\text{rec}}(\mathbf{W}; \Psi)$ for text embeddings. After optimizing the SAE, the encoder extracts sparse
87 latent codes, and the decoder columns provide estimates of monosemantic feature directions [12].

88 **Concept alignment across modalities in sparse autoencoders.** When an SAE is trained on
89 embeddings from both modalities, it is expected to *recover monosemantic features* within each
90 modality and to *align shared concepts* by assigning corresponding image and text features to shared
91 latent coordinates. However, each latent index j is tied to a single column \mathbf{w}_j of the decoder weight
92 matrix in (3). Assigning corresponding features to the same coordinate implicitly assumes that they
93 share a common direction in the joint embedding space. We first examine whether this assumption
94 holds in learned joint embedding spaces.

95 3 Cross-Modal Feature Heterogeneity in Joint Embedding Spaces

96 We examine whether semantically corresponding features share a common direction across image
97 and text modalities in the joint embedding space. When they do not, we refer to this directional
98 mismatch as *cross-modal feature heterogeneity*, which we formally define below.

99 **Definition 1** (Cross-Modal Feature Heterogeneity). *The i -th latent concept exhibits cross-modal*
100 *feature heterogeneity if its image and text feature directions are not perfectly aligned, i.e., $\phi_i \neq \psi_i$.*

101 This definition separates semantic correspondence from directional identity in the joint embedding
102 space. In other words, an image feature and a text feature may represent the same concept while oc-
103 cupying different directions in the joint embedding space. We next examine whether such directional
104 differences appear in joint embedding spaces learned by VLMs.

¹For theory, we use the Top-1 operator, whose behavior can be approximated by Top- K activations with bias terms. In experiments, we use SAEs with Top- K activations for $K > 1$, bias terms, and separate encoder and decoder weights.

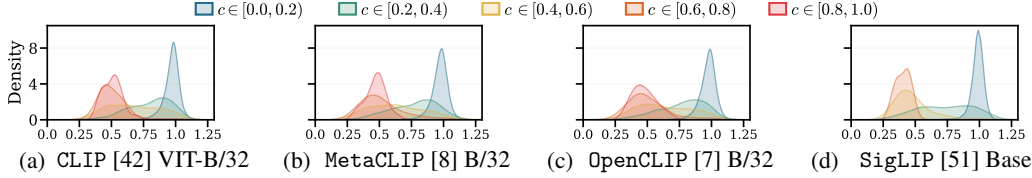


Figure 2: Distribution of cosine distances between image–text feature pairs estimated from embeddings of four VLMs. We group image–text feature pairs by their coactivation correlation $c_{i,j}$ in (5), shown in different colors. The value $c_{i,j}$ measures how strongly the i -th image feature and the j -th text feature coactivate on paired image–text embeddings, so pairs with larger correlations are more likely to represent the same shared concept. Across all models, the distribution does not concentrate near 0 but remains centered around a positive value, even for high-correlation pairs ($c \in [0.8, 1.0]$, shown in blue). This observation supports the presence of cross-modal feature heterogeneity.

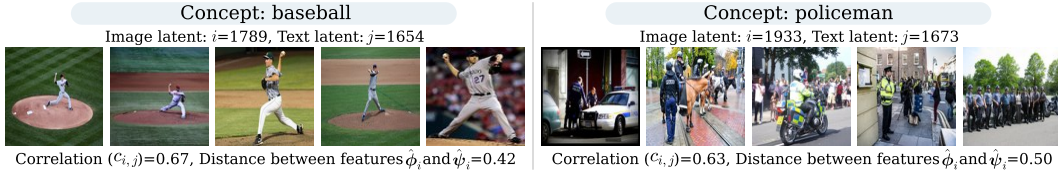


Figure 3: Image examples with high coactivation correlation $c_{i,j}$ from the CLIP setting in Figure 2, where i and j denote the image and text latent indices, respectively. The images reveal coherent concepts such as *baseball* and *policeman*, but the features $\hat{\phi}_i$ and $\hat{\psi}_j$ remain directionally separated.

105 **Measuring differences in feature directions across modalities.** Since the true feature directions
 106 are unobserved, we use decoder columns as their estimates, letting $\hat{\phi}_i$ and $\hat{\psi}_j$ denote the i -th and
 107 j -th decoder columns of SAEs trained on image and text embeddings, respectively.

108 The correspondences between image and text features are also unknown. We therefore use coactivation
 109 correlations on paired image and text embeddings as a proxy for semantic correspondence. Using the
 110 encoders of the SAEs trained on image and text embeddings, we compute latent codes $\tilde{\mathbf{z}}(\mathbf{x})$ and $\tilde{\mathbf{z}}(\mathbf{y})$
 111 for paired embeddings (\mathbf{x}, \mathbf{y}) and form their correlation matrix over the training set:

$$\mathbf{C} := \text{Corr}(\tilde{\mathbf{z}}(\mathbf{x}), \tilde{\mathbf{z}}(\mathbf{y})) \in \mathbb{R}^{m \times m}, \quad c_{i,j} := [\mathbf{C}]_{i,j}. \quad (5)$$

112 Each entry $c_{i,j}$ measures how strongly the i -th image latent and the j -th text latent coactivate on
 113 paired inputs. Larger values are treated as indicating more likely semantic correspondence. For each
 114 pair (i, j) , we measure the cosine distance between the estimated feature directions $\hat{\phi}_i$ and $\hat{\psi}_j$.

115 Figure 2 reports the distribution of cosine distances between estimated image and text feature direc-
 116 tions, grouped by coactivation correlation. The results are computed on MS-COCO [29] embeddings
 117 extracted from four VLMs, including CLIP [42], MetaCLIP [8], OpenCLIP [7], and SigLIP2 [51].
 118 Feature pairs with larger coactivation correlation tend to have smaller distances, but they do not
 119 concentrate near zero distance. This provides empirical support for cross-modal feature heterogeneity,
 120 where semantically corresponding image and text features need not share identical directions in the
 121 joint embedding space. See Appendix E.1 for experimental details.

122 Figure 3 shows examples where matched image and text latents strongly coactivate and respond to the
 123 same concept, such as *baseball* or *policeman*. For example, the *baseball* concept activates the image
 124 latent (with index $i = 1789$) and the text latent (with index $j = 1654$), yielding a high coactivation
 125 correlation, $c_{i,j} = 0.67$. But the cosine distance between the estimated features $\hat{\phi}_i$ and $\hat{\psi}_j$ remains
 126 large at 0.42. These examples show that semantic correspondence need not imply directional identity.

127 Appendix C discusses possible sources of this heterogeneity and relates it to the modality gap,
 128 showing that a nontrivial modality gap implies cross-modal feature heterogeneity.

129 4 Analysis of Sparse Autoencoders under Cross-Modal Feature Heterogeneity

130 We analyze whether an SAE can reconstruct image and text feature directions while assigning
 131 corresponding concepts across modalities to shared latent coordinates. These goals can conflict under
 132 cross-modal feature heterogeneity. Corollary 1 explains why reconstruction induces modality split,

133 Theorem 1 studies the limited capacity case, and Proposition 1 shows why existing alignment methods
 134 reduce modality split by sacrificing reconstruction quality. Proofs are provided in Appendix D.

135 **Why reconstruction induces modality split.** We consider an SAE trained with the reconstruction
 136 loss (4) on both image and text embeddings. Since no alignment term is used, the objective only
 137 encourages accurate reconstruction.

138 **Corollary 1.** *Suppose an SAE has $m \geq 2n$ latent coordinates and is trained by minimizing*
 139 $\mathcal{L}_{\text{rec}}(\mathbf{W}; \Phi) + \mathcal{L}_{\text{rec}}(\mathbf{W}; \Psi)$. *As the sparse factor s in (2) approaches 1, for any permutation*
 140 *matrix $\mathbf{P} \in \mathbb{R}^{m \times m}$, the weight matrix $\hat{\mathbf{W}} := [\Phi \ \Psi \ \mathbf{0}_{d \times (m-2n)}] \mathbf{P}$ is a global minimizer, and*
 141 $\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \Phi) + \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \Psi) = 0$.

142 Corollary 1 shows that reconstruction is minimized by assigning a separate column of the decoder
 143 weight matrix \mathbf{W} to each feature direction. Thus, if the image feature ϕ_i and the text feature ψ_i
 144 represent the same concept but point in different directions, they are assigned to different columns,
 145 say $\hat{\mathbf{w}}_a$ and $\hat{\mathbf{w}}_b$, and activate different latent coordinates, \tilde{z}_a and \tilde{z}_b . This explains why modality split
 146 can be reconstruction optimal rather than a failure of concept recovery. The SAE preserves both
 147 feature directions, but the corresponding image and text concepts appear at different latent indices.

148 **Why limited capacity collapses distinct features.** Corollary 1 assumes enough latent coordinates
 149 to represent all feature directions separately, namely $m \geq 2n$. In practice, the effective number of
 150 coordinates can be smaller due to dead neurons [50, 18, 32], *i.e.*, $\tilde{z}_j(\mathbf{x}) = 0$ for all \mathbf{x} . We therefore
 151 analyze the regime where the SAE has fewer latent coordinates than feature directions, $m < 2n$.

152 **Theorem 1.** *Suppose the SAE has $m < 2n$ latent coordinates and is trained by minimizing*
 153 $\mathcal{L}_{\text{rec}}(\mathbf{W}; \Phi) + \mathcal{L}_{\text{rec}}(\mathbf{W}; \Psi)$. *Define $\mathbf{M}_i := \mathbb{E}[z_i^2 \mid z_i \neq 0] \phi_i \phi_i^\top \in \mathbb{R}^{d \times d}$ and $\mathbf{M}_{n+i} := \mathbb{E}[z_i^2 \mid$
 154 $z_i \neq 0] \psi_i \psi_i^\top \in \mathbb{R}^{d \times d}$ for $i \in [n]$. Let $(\mathbb{A}_1, \dots, \mathbb{A}_m)$ be a partition of $[2n]$ that maximizes*
 155 $\sum_{j \in [m]} \lambda_{\max} \left(\sum_{i \in \mathbb{A}_j} \mathbf{M}_i \right)$, *where $\lambda_{\max}(\cdot)$ denotes the largest eigenvalue of a matrix. As the sparse*
 156 *factor s in (2) approaches 1, a global minimizer $\hat{\mathbf{W}} := [\hat{\mathbf{W}}_1 \ \dots \ \hat{\mathbf{W}}_m]$ is obtained by taking $\hat{\mathbf{W}}_j$ to be*
 157 *a top eigenvector with unit norm of $\sum_{i \in \mathbb{A}_j} \mathbf{M}_i$ for each $j \in [m]$, and $\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \Phi) + \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \Psi) =$
 158 $s^{n-1}(1-s) \left(2 \sum_{i \in [n]} \mathbb{E}[z_i^2 \mid z_i \neq 0] - \sum_{j \in [m]} \lambda_{\max} \left(\sum_{i \in \mathbb{A}_j} \mathbf{M}_i \right) \right) + o(1-s)$.*

159 Theorem 1 shows that when an SAE lacks enough effective coordinates, it must reuse columns of the
 160 decoder weight matrix \mathbf{W} to represent multiple feature directions. In this case, feature directions are
 161 grouped by their geometry in the embedding space rather than by semantic correspondence. Hence,
 162 close but distinct directions may share a column, even when they correspond to different concepts.

163 Under this limited capacity regime, when ϕ_i and ψ_i are close but not identical, the two feature
 164 directions may be grouped together and represented by a single column. This can reduce modality
 165 split by assigning them to a shared coordinate. However, it comes at the cost of worse reconstruction
 166 and poorer feature recovery, because a single column must approximate two distinct directions.

167 **Why alignment during training sacrifices reconstruction.** The previous results show that re-
 168 construction favors separate columns, while alignment favors shared coordinates. We examine this
 169 tension in existing methods that add auxiliary alignment losses during training [14, 25].

170 For clarity, we state a simplified result. In this result, $\lambda^*(\rho)$ and $\lambda^\dagger(\rho)$ are critical alignment strengths
 171 at which the global minimizer changes. Their exact definitions, along with complete statements for
 172 the group-sparse loss [25] and the Iso-Energy alignment loss [14], are given in Propositions 3 and 4.

173 **Proposition 1.** *Consider $n = 1$ and define $\rho := \phi^\top \psi \in (0, 1)$. Suppose the SAE has $m = 2$*
 174 *latent coordinates, and each column of \mathbf{W} has unit norm. For the group-sparse loss [25], where*
 175 *λ is the weight of the auxiliary alignment term, a global minimizer is $[\phi \ \psi]$ when $\lambda \in [0, \lambda^*(\rho))$,*
 176 $\left[\frac{\phi + \psi}{\|\phi + \psi\|} \ \mathbf{0}_d \right]$ *when $\lambda \in (\lambda^*(\rho), \lambda^\dagger(\rho))$, and satisfies $\sigma(\mathbf{W}^\top \phi) = \sigma(\mathbf{W}^\top \psi) = \mathbf{0}$ when $\lambda \in$*
 177 $(\lambda^\dagger(\rho), \infty)$. *The reconstruction losses for each λ regime are 0, $(1-\rho)\mathbb{E}[z^2]$, and $2\mathbb{E}[z^2]$, respectively.*

178 Proposition 1 illustrates the trade-off between reconstruction and alignment. Small λ keeps two
 179 separate columns, $[\phi \ \psi]$, yielding zero reconstruction loss but preserving modality split. Larger
 180 λ collapses the two directions into a single shared column of the decoder weight matrix, $\frac{\phi + \psi}{\|\phi + \psi\|}$,

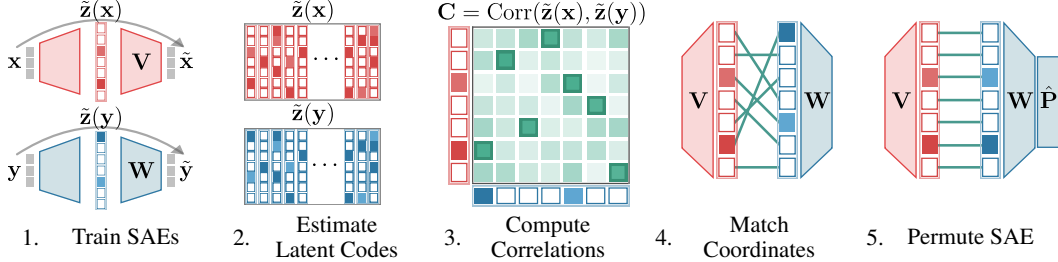


Figure 4: Overview of our approach. We first train modality-specific SAEs to reconstruct image and text embeddings, obtaining latent codes $\tilde{z}(x)$ and $\tilde{z}(y)$. We then compute their correlation matrix C in (5) over paired training data and apply the Hungarian algorithm to obtain the assignment \hat{P} in (6). Finally, we reindex the text latent coordinates using \hat{P} , aligning corresponding image and text features to shared latent indices without altering the learned feature directions.

181 improving alignment but incurring nonzero reconstruction loss. For the group-sparse loss, an even
 182 larger λ leads to a degenerate solution where both features become inactive.

183 5 Modality-Specific Sparse Autoencoders and Post-Hoc Alignment

184 Motivated by the analysis in Section 4, we propose a two-stage method. First, we train separate SAEs
 185 for image and text embeddings to preserve modality-specific feature directions. Then, we align their
 186 latent coordinates using activation correlations. Figure 4 provides an overview of this procedure.

187 **Modality-specific sparse autoencoders.** We train separate SAEs for image and text embeddings,
 188 with weight matrices V and W , respectively, using the reconstruction loss in (4). We minimize
 189 $\mathcal{L}_{\text{rec}}(V; \Phi)$ and $\mathcal{L}_{\text{rec}}(W; \Psi)$ separately, without any cross-modal alignment constraint. This prevents
 190 image and text feature directions from competing for the same columns of the weight matrix, allowing
 191 each modality to preserve its own directions before cross-modal correspondences are identified.

192 **Post-hoc alignment.** Because the latent coordinates of the two SAEs are not naturally aligned, we
 193 use the correlation matrix C in (5) to identify corresponding image and text latents. Specifically, we
 194 find a permutation matrix \hat{P} that maximizes the total correlation between aligned coordinates,

$$\hat{P} \in \arg \max_{P \in \mathcal{P}_m} \text{tr}(CP), \quad (6)$$

195 where \mathcal{P}_m denotes the set of all $m \times m$ permutation matrices. The trace sums the correlations
 196 of image and text coordinates matched by P . We compute this assignment using the Hungarian
 197 algorithm [26]². We then apply the permutation only to the SAE for the text modality as

$$\tilde{z}(y; W) := \sigma(\hat{P}^\top W^\top y), \quad \tilde{y}(y; W) := W\hat{P}\tilde{z}(y; W). \quad (7)$$

198 This alignment simply permutes the learned features of the text modality without changing their direc-
 199 tions. Consequently, it preserves the feature geometry of each modality while assigning corresponding
 200 features from the image and text modalities to the same coordinates in the latent codes.

201 **Inference.** Embeddings are encoded into latent codes using modality-specific SAEs, and the
 202 permutation is applied to align their indices as in (7). This alignment enables cross-modal tasks such
 203 as retrieval in the aligned latent space, while encoding and decoding remain modality-specific.

204 6 Experiments

205 This section validates the theoretical findings in Section 4 and evaluates the approach proposed in
 206 Section 5. Section 6.1 uses synthetic embeddings where ground-truth feature directions and cross-
 207 modal correspondences are known, allowing us to directly test the theoretical predictions. Section 6.2
 208 evaluates our method on real-world datasets to assess its effectiveness on cross-modal tasks.

209 The codes used in our experiments are available in the anonymous repository.

²Dead latent coordinates are excluded in practice. Thus, (6) is applied to the corresponding submatrix of C .



Figure 5: Comparison between a shared SAE trained on both modalities (green) and modality-specific SAEs trained separately for each modality (orange) as the cross-modal feature distance varies. Each point corresponds to SAEs trained on synthetic embeddings sampled from ground-truth feature matrices Φ and Ψ , where the x-axis shows the cosine distance between corresponding image and text feature directions. When corresponding directions are distinct but geometrically close, the shared SAE tends to collapse them into a single learned direction, leading to higher reconstruction and feature recovery errors. In contrast, modality-specific SAEs avoid this collapse and better preserve feature directions under cross-modal feature heterogeneity, which corresponds to nonzero distance.

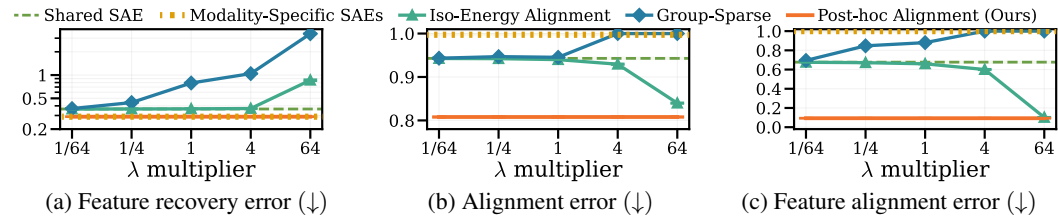


Figure 6: Comparison between Post-hoc Alignment method (orange) and two auxiliary-loss baselines: Iso-Energy alignment [14] (teal) and group-sparse loss [25] (blue). The coefficient of auxiliary-loss λ is shown as a multiplier of the default value in the original papers. For reference, we also include the shared SAE (green) and modality-specific SAEs (yellow), neither of which depend on λ . Our method achieves both lower feature recovery error and lower alignment error than the baselines.

210 6.1 Validation on Synthetic Embeddings

211 We validate two key design choices of the method proposed in Section 5. First, we test whether
 212 training modality-specific SAEs better preserves image and text feature directions than training a
 213 shared SAE. Second, we test whether post-hoc alignment can align corresponding latent coordinates
 214 without sacrificing feature recovery. These design choices are motivated by the analysis in Section 4.

215 **Setup.** To test the first choice, we compare a *shared SAE*, which uses one weight matrix for both
 216 modalities, with *modality-specific SAEs*, which use separate weight matrices for image and text
 217 embeddings. For a fair comparison, we match the total number of trainable parameters by using
 218 $m/2$ latent coordinates for each modality-specific SAE when the shared SAE uses m . Thus, any
 219 performance difference reflects how capacity is distributed across modalities, rather than model size.

220 To test the second choice, we compare our post-hoc alignment method with two auxiliary-loss
 221 baselines [14, 25]. We report the average over three independent runs. Details on synthetic embedding
 222 generation and training setup are provided in Appendix E.2.

223 **Metrics.** We evaluate each method along three axes: *Reconstruction Error* and *Feature Recovery*
 224 *Error* for reconstruction quality; *Alignment Error* and *Feature Alignment Error* for cross-modal
 225 alignment; and *Feature Collapse Rate* for whether same-concept features across modalities are
 226 collapsed into a single learned direction. See Appendix E.3 for formal definitions.

227 Modality-specific sparse autoencoders better preserve feature directions under heterogeneity.

228 Figure 5 reports the results obtained by varying the cosine distance between corresponding cross-
 229 modal features ϕ_i and ψ_i . As shown in Figure 5a, the shared SAE exhibits feature collapse when
 230 image and text feature directions are geometrically close but not identical. Modality-specific SAEs
 231 avoid this collapse because each modality has its own SAE. Figures 5b and 5c show that this collapse
 232 leads to larger reconstruction and feature recovery errors. This is consistent with the analysis in
 233 Section 4, where representing different feature directions with a single decoder column increases

Table 1: Reconstruction and alignment quality of various methods. Reconstruction quality is measured by mean squared error between the input embedding and its reconstruction. Cross-modal alignment is evaluated by Recall@ k for $k \in \{1, 5, 10\}$ on image-to-text and text-to-image retrieval tasks on MS-COCO [29], and by top-1 zero-shot classification accuracy on ImageNet1K [13]. The best and second-best results for each metric are shown in bold and underlined, respectively.

Methods	MS-COCO [29]							ImageNet1K [13]	
	Recon. (\downarrow)	Image-to-Text (\uparrow)			Text-to-Image (\uparrow)			Recon. (\downarrow)	Zero-shot (\uparrow)
	MSE	Recall@1	Recall@5	Recall@10	Recall@1	Recall@5	Recall@10	MSE	Accuracy
Shared SAE	0.090	6.1 (± 0.3)	13.3 (± 0.8)	17.8 (± 0.8)	3.4 (± 0.2)	9.5 (± 0.5)	14.0 (± 0.8)	0.118	15.7 (± 1.3)
+ Iso-Energy alignment loss	0.091	4.7 (± 1.9)	10.6 (± 3.4)	14.6 (± 4.4)	2.7 (± 0.7)	7.4 (± 2.1)	10.8 (± 3.0)	0.118	13.0 (± 2.3)
+ Group-sparse loss	0.105	<u>7.1</u> (± 0.1)	<u>16.7</u> (± 0.4)	<u>23.8</u> (± 0.3)	<u>4.3</u> (± 0.1)	<u>12.2</u> (± 0.1)	<u>18.3</u> (± 0.2)	0.134	26.6 (± 0.3)
Modality-Specific SAEs	0.089	0.0 (± 0.0)	0.1 (± 0.1)	0.2 (± 0.1)	0.0 (± 0.0)	0.1 (± 0.0)	0.1 (± 0.1)	0.116	0.1 (± 0.0)
+ Post-hoc Alignment (Ours)	0.089	16.0 (± 0.5)	34.1 (± 1.4)	44.5 (± 1.6)	11.4 (± 0.4)	27.2 (± 1.2)	37.0 (± 1.3)	0.116	<u>25.1</u> (± 0.3)

reconstruction loss. At zero cosine distance, the image and text feature directions coincide, so there is no cross-modal feature heterogeneity. The shared SAE can therefore represent the same concept with a single direction, which explains its lower reconstruction and feature recovery errors in this idealized case. Overall, these results show that modality-specific SAEs better preserve feature directions when semantically corresponding features differ across modalities, leading to better feature recovery. Appendix F.2 further shows that this phenomenon also appears in SAEs trained on VLM embeddings.

Post-hoc alignment improves alignment without sacrificing feature recovery. Figure 6 reports the results obtained by varying the auxiliary-loss strength λ . We fix the cosine distance between corresponding cross-modal features at 0.5 to reflect the level of cross-modal feature heterogeneity observed in Figure 2. For the Iso-Energy alignment loss [14], increasing λ improves cross-modal alignment but degrades feature recovery, consistent with the reconstruction-alignment trade-off shown in Proposition 4. For the group-sparse loss [25], increasing λ eventually pushes the SAE into a degenerate regime where both feature recovery and alignment deteriorate, consistent with Proposition 3. In contrast, our post-hoc approach preserves the feature recovery achieved by modality-specific SAEs before alignment, while achieving better alignment than the auxiliary-loss baselines. These results support decoupling cross-modal alignment from the reconstruction objective.

6.2 Evaluation on Real-World Data

We evaluate our method on real-world image-text embeddings extracted from a pre-trained VLM. We assess whether latent codes support cross-modal tasks and preserve semantically coherent activations. In particular, following prior work [25], we evaluate whether estimated latent codes \tilde{z} in (5) provide useful representations for cross-modal retrieval, concept steering, and monosemanticity evaluation.

Experimental protocol. We use CC-3M [48] as a paired image-text dataset. From this dataset, we extract image and text embeddings using the CLIP model [42] with the ViT-B/32 architecture, and train SAEs on the resulting paired embeddings. We compare the method proposed in Section 5 with two auxiliary-loss baselines [14, 25]. To ensure a fair comparison between modality-specific and shared SAE approaches, all methods use the same total number of SAE parameters. We follow prior training protocols [38] and report averages over three runs. Details are provided in Appendix E.4.

Our method improves cross-modal performance while preserving reconstruction. We evaluate reconstruction quality and performance across two cross-modal tasks. Reconstruction quality is measured by the mean squared error between input embeddings and their reconstructions. For cross-modal retrieval on MS-COCO, we rank image and text latent codes by cosine similarity and report Recall@ k . For zero-shot image classification on ImageNet1K [13], we report top-1 accuracy.

Table 1 reports the results for all methods. Our method achieves the lowest reconstruction error and the strongest cross-modal retrieval performance, while remaining competitive on zero-shot classification. In cross-modal retrieval, our method outperforms the strongest baseline by 8.9 points in image-to-text Recall@1 (16.0 vs. 7.1) and by 7.1 points in text-to-image Recall@1 (11.4 vs. 4.3). These results show that post-hoc alignment improves cross-modal performance without sacrificing the reconstruction quality of modality-specific SAEs. Additional results across a broader range of VLM sizes and families are provided in Appendix F.

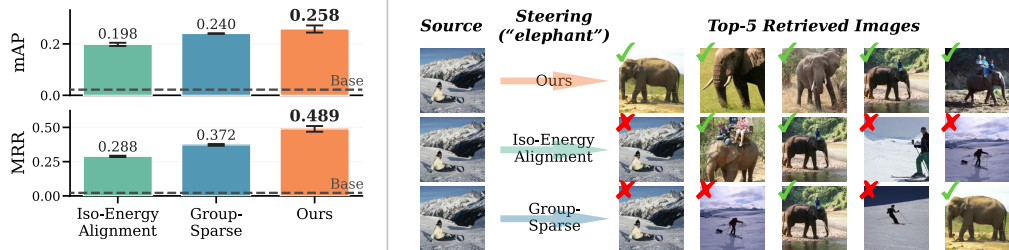


Figure 7: Image retrieval under concept latent steering on MS-COCO. We steer source image embeddings using target feature directions identified through aligned latent coordinates. Retrieval performance is measured by mean average precision (mAP) and mean reciprocal rank (MRR) over retrieved images containing the target concept. Qualitative examples show source images and their top retrieved images after steering. Our method steers images toward target concepts more effectively.

273 **Our method enables more effective cross-modal concept steering.** We next evaluate whether the
 274 aligned latent coordinates support controllable steering of object concepts. We treat each of the 80
 275 object categories in COCO as a target concept. For each target concept, we compute the mean activation
 276 of each text latent coordinate on captions that mention the concept and on randomly sampled captions
 277 that do not. We then select the text latent coordinate with the largest difference between the two mean
 278 activations, which identifies the coordinate most associated with the target concept. Using the aligned
 279 image coordinate, we use the corresponding decoder column in the image SAE as the steering vector.

280 For each target concept, we select 100 source images from the test set that do not contain the target
 281 category. We add the steering vector to each source image embedding, producing a steered embedding
 282 intended to move it toward the target concept. We then rank all test image embeddings by cosine
 283 similarity to each steered embedding and retrieve the nearest images. Successful steering should
 284 retrieve images that contain the target concept, even though the source images do not.

285 Figure 7 reports the retrieval performance and qualitative examples of concept steering. Our method
 286 achieves the highest mean average precision and mean reciprocal rank over retrieved test images that
 287 contain the target concept. The qualitative examples further show that our method retrieves target-
 288 concept images more consistently. These results indicate that post-hoc latent coordinate alignment
 289 enables more reliable control than enforcing alignment through an auxiliary loss during training.

290 **Our method preserves more semantically coherent latents.** We examine whether the learned latent
 291 codes are activated by semantically coherent inputs. We use the monosemanticity score [37] on the
 292 validation set of CC-3M. This score measures whether inputs that activate the same latent are close to one
 293 another in an external embedding space. We use the MetaCLIP model [8] as the external encoder. Larger
 294 values indicate that the latent responds to a more semantically coherent concept.
 295
 296
 297
 298

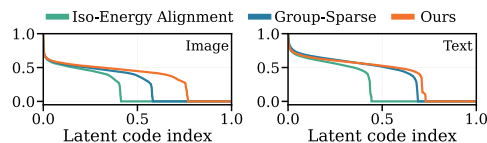


Figure 8: Monosemanticity scores for coordinates sorted in descending order, shown separately for image and text modalities.

299 Figure 8 reports monosemanticity scores for each latent coordinate, sorted in descending order. Our
 300 method maintains higher scores across more latent coordinates in both modalities. This indicates that
 301 preserving modality-specific feature directions gives more semantically coherent latents.

302 7 Conclusion

303 We studied how features are organized across modalities in joint embedding spaces of vision–language
 304 models. We characterized *cross-modal feature heterogeneity*, where the same semantic concept can
 305 have different feature directions across image and text modalities. We showed that this heterogeneity
 306 can give rise to modality split in SAEs, and that existing auxiliary-loss approaches trade reconstruction
 307 quality for latent alignment. Motivated by this limitation, we proposed a simple approach that trains
 308 modality-specific SAEs and aligns their latent codes through coactivation. In experiments, our method
 309 preserves reconstruction quality and improves cross-modal alignment. Overall, effective alignment in
 310 joint embedding spaces should preserve each modality’s feature geometry.

311 Due to space constraints, a discussion of limitations and future directions is provided in Appendix B.

312 References

- 313 [1] Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier
314 probes. In *International Conference on Learning Representations*, 2017.
- 315 [2] Andy Arditì, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and
316 Neel Nanda. Refusal in language models is mediated by a single direction. In *Conference on
317 Neural Information Processing Systems*, 2024.
- 318 [3] Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. Linear algebraic
319 structure of word senses, with applications to polysemy. *Transactions of the Association for
320 Computational Linguistics*, 6:483–495, 07 2018.
- 321 [4] Manuel Brack, Felix Friedrich, Dominik Hintersdorf, Lukas Struppek, Patrick Schramowski, and
322 Kristian Kersting. Sega: Instructing text-to-image models using semantic guidance. *Conference
323 on Neural Information Processing Systems*, 2023.
- 324 [5] Bart Bussmann, Patrick Leask, and Neel Nanda. Batchtopk sparse autoencoders. *arXiv preprint
325 arXiv:2412.06410*, 2024.
- 326 [6] Zixiang Chen, Yihe Deng, Yuanzhi Li, and Quanquan Gu. Understanding transferable repre-
327 sentation learning and zero-shot transfer in CLIP. In *International Conference on Learning
328 Representations*, 2024.
- 329 [7] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade
330 Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling
331 laws for contrastive language-image learning. In *Conference on Computer Vision and Pattern
332 Recognition*, 2023.
- 333 [8] Yung-Sung Chuang, Yang Li, Dong Wang, Ching-Feng Yeh, Kehan Lyu, Ramya Raghavendra,
334 James R. Glass, LIFEI HUANG, Jason E Weston, Luke Zettlemoyer, Xinlei Chen, Zhuang Liu,
335 Saining Xie, Wen tau Yih, Shang-Wen Li, and Hu Xu. Meta CLIP 2: A worldwide scaling
336 recipe. In *Conference on Neural Information Processing Systems*, 2026.
- 337 [9] Bilal Chughtai, Lawrence Chan, and Neel Nanda. A toy model of universality: Reverse
338 engineering how networks learn group operations. In *International Conference on Machine
339 Learning*, 2023.
- 340 [10] Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià
341 Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. In
342 *Conference on Neural Information Processing Systems*, 2023.
- 343 [11] Valérie Costa, Thomas Fel, Ekdeep Singh Lubana, Bahareh Tolooshams, and Demba E. Ba.
344 From flat to hierarchical: Extracting sparse representations with matching pursuit. In *Conference
345 on Neural Information Processing Systems*, 2026.
- 346 [12] Jingyi Cui, Qi Zhang, Yifei Wang, and Yisen Wang. On the limits of sparse autoencoders:
347 A theoretical framework and reweighted remedy. In *International Conference on Learning
348 Representations*, 2026.
- 349 [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
350 hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, 2009.
- 351 [14] Grégoire DHIMOÏLA, Thomas Fel, Victor Boutin, and Agustin Martin Picard. Cross-modal
352 redundancy and the geometry of vision–language embeddings. In *International Conference on
353 Learning Representations*, 2026.
- 354 [15] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,
355 Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep
356 Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt,
357 Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish,
358 and Chris Olah. A mathematical framework for transformer circuits, 2021. URL <https://transformer-circuits.pub/2021/framework/index.html>.
359

- 360 [16] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna
361 Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of
362 superposition. *arXiv preprint*, 2022.
- 363 [17] Sedigheh Eslami and Gerard de Melo. Mitigate the gap: Improving cross-modal alignment in
364 CLIP. In *International Conference on Learning Representations*, 2025.
- 365 [18] Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya
366 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. In *International
367 Conference on Learning Representations*, 2025.
- 368 [19] Ian J Goodfellow, Aaron Courville, and Yoshua Bengio. Large-scale feature learning with spike-
369 and-slab sparse coding. In *International Conference on Machine
370 Learning*, 2012.
- 371 [20] Eleonora Grassucci, Giordano Cicchetti, Emanuele Frasca, Aurelio Uncini, and Danilo Com-
372 miniello. Closing the modality gap aligns group-wise semantics. In *International Conference
373 on Learning Representations*, 2026.
- 374 [21] Ruben Härle, Felix Friedrich, Manuel Brack, Stephan Wäldchen, Björn Deiseroth, Patrick
375 Schramowski, and Kristian Kersting. Measuring and guiding monosemanticity. In *Conference
376 on Neural Information Processing Systems*, 2025.
- 377 [22] John Hewitt and Christopher D Manning. A structural probe for finding syntax in word
378 representations. In *North American Chapter of the Association for Computational Linguistics*,
379 2019.
- 380 [23] Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. Sparse
381 autoencoders find highly interpretable features in language models. In *International Conference
382 on Learning Representations*, 2023.
- 383 [24] Qian Jiang, Changyou Chen, Han Zhao, Liqun Chen, Qing Ping, Son Dinh Tran, Yi Xu, Belinda
384 Zeng, and Trishul Chilimbi. Understanding and constructing latent modality structures in multi-
385 modal representation learning. In *Conference on Computer Vision and Pattern Recognition*,
386 2023.
- 387 [25] Chiraag Kaushik, Davis Barch, and Andrea Fanelli. Decomposing multimodal embedding spaces
388 with group-sparse autoencoders. In *International Conference on Learning Representations*,
389 2026.
- 390 [26] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics
391 quarterly*, 2(1-2):83–97, 1955.
- 392 [27] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-
393 time intervention: Eliciting truthful answers from a language model. In *Conference on Neural
394 Information Processing Systems*, 2023.
- 395 [28] Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Zou. Mind the
396 gap: Understanding the modality gap in multi-modal contrastive representation learning. In
397 *Conference on Neural Information Processing Systems*, 2022.
- 398 [29] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
399 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European
400 Conference on Computer Vision*, 2014.
- 401 [30] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In
402 *Conference on Neural Information Processing Systems*, 2023.
- 403 [31] Sheng Liu, Haotian Ye, and James Zou. Reducing hallucinations in large vision-language
404 models via latent space steering. In *International Conference on Learning Representations*,
405 2025.
- 406 [32] Lu Lu, Yeonjong Shin, Yanhui Su, and George Em Karniadakis. Dying relu and initialization:
407 Theory and numerical examples. *Communications in Computational Physics*, 28(5):1671–1706,
408 2020.

- 409 [33] Alireza Makhzani and Brendan Frey. K-sparse autoencoders. In *International Conference on*
410 *Learning Representations*, 2014.
- 411 [34] Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller.
412 Sparse feature circuits: Discovering and editing interpretable causal graphs in language models.
413 In *International Conference on Learning Representations*, 2025.
- 414 [35] Marco Mistretta, Alberto Baldrati, Lorenzo Agnolucci, Marco Bertini, and Andrew D. Bagdanov.
415 Cross the gap: Exposing the intra-modal misalignment in CLIP via modality inversion. In
416 *International Conference on Learning Representations*, 2025.
- 417 [36] Tuomas Oikarinen and Tsui-Wei Weng. CLIP-dissect: Automatic description of neuron repre-
418 sentations in deep vision networks. In *International Conference on Learning Representations*,
419 2023.
- 420 [37] Mateusz Pach, Shyamgopal Karthik, Quentin Bouniot, Serge Belongie, and Zeynep Akata.
421 Sparse autoencoders learn monosemantic features in vision-language models. In *Conference on*
422 *Neural Information Processing Systems*, 2026.
- 423 [38] Isabel Papadimitriou, Huangyuan Su, Thomas Fel, Sham M. Kakade, and Stephanie Gil.
424 Interpreting the linear structure of vision-language model embedding spaces. In *Conference on*
425 *Language Modeling*, 2025.
- 426 [39] Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the
427 geometry of large language models. In *International Conference on Machine Learning*, 2024.
- 428 [40] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip:
429 Text-driven manipulation of stylegan imagery. In *International Conference on Computer Vision*,
430 2021.
- 431 [41] Qi Qian, Yuanhong Xu, and Juhua Hu. Intra-modal proxy learning for zero-shot visual catego-
432 rization with clip. In *Conference on Neural Information Processing Systems*, 2023.
- 433 [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-
434 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
435 Sutskever. Learning transferable visual models from natural language supervision. In *Internati-
436 onal Conference on Machine Learning*, 2021.
- 437 [43] Sameera Ramasinghe, Violetta Shevchenko, Gil Avraham, and Ajanthan Thalaiyasingam.
438 Accept the modality gap: An exploration in the hyperbolic space. In *Conference on Computer*
439 *Vision and Pattern Recognition*, 2024.
- 440 [44] Sukrut Rao, Sweta Mahajan, Moritz Böhle, and Bernt Schiele. Discover-then-name: Task-
441 agnostic concept bottlenecks via automated concept discovery. In *European Conference on*
442 *Computer Vision*, 2024.
- 443 [45] Nina Rimskey, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner.
444 Steering llama 2 via contrastive activation addition. In *Annual Meeting of the Association for*
445 *Computational Linguistics*, 2024.
- 446 [46] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
447 resolution image synthesis with latent diffusion models. In *Conference on Computer Vision and*
448 *Pattern Recognition*, 2022.
- 449 [47] Simon Schrodi, David T. Hoffmann, Max Argus, Volker Fischer, and Thomas Brox. Two
450 effects, one trigger: On the modality gap, object bias, and information imbalance in contrastive
451 vision-language models. In *International Conference on Learning Representations*, 2025.
- 452 [48] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A
453 cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Association for*
454 *Computational Linguistics*, 2018.
- 455 [49] Abdul-Saboor Sheikh, Jacquelyn A Shelton, and Jörg Lücke. A truncated em approach for
456 spike-and-slab sparse coding. *Journal of Machine Learning Research*, 15(1):2653–2687, 2014.

- 457 [50] Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen,
458 Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L.
459 Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Ed-
460 ward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling
461 monosemanticity: Extracting interpretable features from claude 3 sonnet, 2024. URL <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>.
462
- 463 [51] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alab-
464 dulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip
465 2: Multilingual vision-language encoders with improved semantic understanding, localization,
466 and dense features. *arXiv preprint*, 2025.
- 467 [52] Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt.
468 Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In
469 *International Conference on Learning Representations*, 2023.
- 470 [53] Shin’ya Yamaguchi, Dewei Feng, Sekitoshi Kanai, Kazuki Adachi, and Daiki Chijiwa. Post-
471 pre-training for modality alignment in vision-language foundation models. In *Conference on*
472 *Computer Vision and Pattern Recognition*, 2025.
- 473 [54] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui
474 Wu. Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine*
475 *Learning Research*, 2022.
- 476 [55] Vladimir Zai grajew, Hubert Baniecki, and Przemyslaw Biecek. Interpreting CLIP with hierar-
477 chical sparse autoencoders. In *International Conference on Machine Learning*, 2025.
- 478 [56] Kaichen Zhang, Yifei Shen, Bo Li, and Ziwei Liu. Large multi-modal models can interpret
479 features in large multi-modal models. In *International Conference on Computer Vision*, 2025.
- 480 [57] Yuhui Zhang, Elaine Sui, and Serena Yeung. Connect, collapse, corrupt: Learning cross-modal
481 tasks with uni-modal data. In *International Conference on Learning Representations*, 2024.

482 A Related Work

483 **Mechanistic interpretability and linear representation hypothesis.** Mechanistic interpretability
484 aims to reverse-engineer neural networks, often by interpreting their learned representations as
485 human-understandable mechanisms [10, 15, 52, 9]. A key assumption underlying this line of work is
486 the *linear representation hypothesis* [16, 39], which posits that each semantically coherent concept
487 is encoded as a unique direction in representation space, and that embeddings can be expressed as
488 linear combinations of such directions. This hypothesis is supported by diverse empirical evidence
489 from linear probing [1, 22] and steering interventions [27, 2, 45]. Recent work has begun extending
490 this hypothesis to multimodal representations [40, 4, 31], opening avenues for understanding feature
491 structure across different modalities.

492 **Sparse autoencoders for interpreting representations.** SAEs have shown great promise for
493 extracting monosemantic features from embeddings of language models [23, 50, 18, 34, 21]. Building
494 on these advances, SAEs have been further applied to VLMs [44, 55, 37], especially joint embedding
495 models trained with contrastive objectives [42]. Recent work has further examined how these features
496 align or separate across modalities [11, 55, 56, 38, 37], motivating an analysis of their underlying
497 structural properties in joint embedding spaces.

498 **Joint representations of vision-language models.** A prominent phenomenon in joint VLM rep-
499 resentations is the *modality gap*, where image and text embeddings occupy disjoint regions of the
500 shared space [28]. Prior work has analyzed this gap from multiple geometric, distributional, and
501 optimization-level perspectives [28, 57, 47]. One line of work argues that the modality gap partly
502 reflects modality-specific features in multimodal representations, rather than merely a failure of
503 alignment [24, 43, 41]. Another line focuses on cross-modal misalignment of shared features and
504 studies how reducing such mismatch improves downstream task performance [17, 53, 20, 35].

505 Recent SAE-based analyses further reveal *modality split* [38], where the same concept activates
506 different latent codes across modalities. Existing remedies force shared latent activations [25, 14],
507 implicitly assuming feature directions are aligned (*i.e.*, $\phi_i = \psi_i$) and overlooking directional
508 misalignment. When shared concepts are encoded along distinct directions, forcing a single shared
509 SAE inevitably degrades reconstruction. We instead preserve modality-specific directions by using
510 modality-specific SAEs and align the corresponding features post hoc.

511 B Limitations and Future Work

512 Our analysis relies on the linear representation hypothesis and simplified generative assumptions,
513 which may not fully capture the complexity of real-world VLM embeddings. While our post-hoc
514 alignment improves empirical performance, it assumes that feature correspondence can be inferred
515 from coactivation statistics, which may become less accurate in noisy or weakly aligned settings.

516 Future work includes extending the analysis beyond the linear representation hypothesis and studying
517 alignment under more complex forms of heterogeneity. While our current experiments focus on
518 CLIP-like contrastive VLMs, another important direction is to evaluate whether the proposed method
519 remains effective for representations extracted from broader VLM architectures, including LLM-based
520 autoregressive VLMs such as LLaVA. Applying these ideas to larger-scale VLMs and downstream
521 tasks may further clarify the role of feature-level alignment in multimodal systems. Moreover,
522 post-hoc alignment could be used as an initialization for subsequent fine-tuning.

523 C Additional Discussion on Cross-Modal Feature Heterogeneity

524 **A possible source of cross-modal feature heterogeneity.** We next discuss why cross-modal feature
 525 directions may be unevenly aligned in VLMs trained with objectives such as contrastive learning.
 526 Such objectives encourage cross-modal similarity for paired samples, typically by increasing $\mathbb{E}[\mathbf{x}^\top \mathbf{y}]$.
 527 As the sparsity level $s \rightarrow 1$, we have

$$\mathbb{E}[\mathbf{x}^\top \mathbf{y}] = \text{tr}(\mathbf{\Phi}^\top \mathbf{\Psi} \mathbb{E}[\mathbf{z}\mathbf{z}^\top]) = s^{n-1}(1-s) \sum_{i \in [n]} \phi_i^\top \psi_i \cdot \mathbb{E}[z_i^2 | z_i \neq 0] + o(1-s).$$

528 This expression suggests that each feature pair contributes to positive similarity in proportion to the
 529 second moment of its latent coordinate. Consequently, under finite data or finite model capacity,
 530 frequently activated features may receive stronger alignment pressure, whereas rarer features may
 531 remain less aligned. This provides one possible source of cross-modal feature heterogeneity.

532 **Relation to modality gap.** The following proposition shows that cross-modal feature heterogeneity
 533 is unavoidable whenever the embedding pair exhibits a nontrivial modality gap.

534 **Proposition 2.** *If there exists a modality gap between an embedding pair, i.e., $\mathbb{E}[\cos(\mathbf{x}, \mathbf{y})] < 1$, then*
 535 *there exists at least one latent concept that exhibits cross-modal feature heterogeneity.*

536 *Proof.* For contradiction, assume that there is no cross-modal feature heterogeneity. Then $\phi_i = \psi_i$
 537 for all $i \in [n]$, which implies $\mathbf{\Phi} = \mathbf{\Psi}$. Hence, the embeddings satisfy $\mathbf{x} = \mathbf{\Phi}\mathbf{z} = \mathbf{\Psi}\mathbf{z} = \mathbf{y}$ almost
 538 surely. This gives $\cos(\mathbf{x}, \mathbf{y}) = 1$ almost surely, contradicting the condition that $\mathbb{E}[\cos(\mathbf{x}, \mathbf{y})] < 1$.
 539 Therefore, there must exist at least one $i \in [n]$ such that $\phi_i \neq \psi_i$. \square

540 This result formalizes that a modality gap at the embedding must be reflected by at least one directional
 541 mismatch between corresponding image and text features under the linear representation hypothesis.
 542 The main question studied in this paper is how such heterogeneity affects multimodal SAEs.

543 D Proofs of Theoretical Results

544 We begin by clarifying notation. We use $[n] := \{1, 2, \dots, n\}$. Let $\mathbf{0}_d$ and $\mathbf{0}_{m \times n}$ denote the zero
 545 vector in \mathbb{R}^d and the $m \times n$ zero matrix, respectively. For a matrix \mathbf{Z} , $[\mathbf{Z}]_{[:,i]}$ denotes its i -th column.

546 **Theorem 1.** *Suppose the SAE has $m < 2n$ latent coordinates and is trained by minimizing*
 547 $\mathcal{L}_{\text{rec}}(\mathbf{W}; \mathbf{\Phi}) + \mathcal{L}_{\text{rec}}(\mathbf{W}; \mathbf{\Psi})$. *Define $\mathbf{M}_i := \mathbb{E}[z_i^2 | z_i \neq 0] \phi_i \phi_i^\top \in \mathbb{R}^{d \times d}$ and $\mathbf{M}_{n+i} := \mathbb{E}[z_i^2 |$
 548 $z_i \neq 0] \psi_i \psi_i^\top \in \mathbb{R}^{d \times d}$ for $i \in [n]$. Let $(\mathbb{A}_1, \dots, \mathbb{A}_m)$ be a partition of $[2n]$ that maximizes
 549 $\sum_{j \in [m]} \lambda_{\max}(\sum_{i \in \mathbb{A}_j} \mathbf{M}_i)$, where $\lambda_{\max}(\cdot)$ denotes the largest eigenvalue of a matrix. As the sparse
 550 factor s in (2) approaches 1, a global minimizer $\hat{\mathbf{W}} := [\hat{\mathbf{W}}_1 \dots \hat{\mathbf{W}}_m]$ is obtained by taking $\hat{\mathbf{W}}_j$ to be
 551 a top eigenvector with unit norm of $\sum_{i \in \mathbb{A}_j} \mathbf{M}_i$ for each $j \in [m]$, and $\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \mathbf{\Phi}) + \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \mathbf{\Psi}) =$
 552 $s^{n-1}(1-s) \left(2 \sum_{i \in [n]} \mathbb{E}[z_i^2 | z_i \neq 0] - \sum_{j \in [m]} \lambda_{\max}(\sum_{i \in \mathbb{A}_j} \mathbf{M}_i) \right) + o(1-s)$.*

553 *Proof.* Let $\theta_i := \phi_i$ and $\theta_{n+i} := \psi_i$ for $i \in [n]$, and $\mu'_i := \mathbb{E}[z_i^2 | z_i \neq 0] > 0$, with $\mu'_{n+i} := \mu'_i$. By
 554 the same expansion as in (9), in the sparsity regime $s \rightarrow 1$,

$$\mathcal{L}_{\text{rec}}(\mathbf{W}; \mathbf{\Phi}) + \mathcal{L}_{\text{rec}}(\mathbf{W}; \mathbf{\Psi}) = s^{n-1}(1-s) \sum_{i \in [2n]} \mu'_i \|\theta_i - \mathbf{W}\sigma(\mathbf{W}^\top \theta_i)\|_2^2 + o(1-s). \quad (8)$$

555 For each $i \in [2n]$, let $\hat{\mathbf{W}}_i \in \arg \max_{\mathbf{W} \in \{\mathbf{W}_j\}_{j \in [m]}} \mathbf{W}^\top \theta_i$. Then

$$\sum_{i \in [2n]} \mu'_i \|\theta_i - \mathbf{W}\sigma(\mathbf{W}^\top \theta_i)\|_2^2 = \sum_{i \in [2n]} \mu'_i \|\theta_i - \hat{\mathbf{W}}_i \hat{\mathbf{W}}_i^\top \theta_i\|_2^2 = \sum_{i \in [2n]} \mu'_i (1 - (\theta_i^\top \hat{\mathbf{W}}_i)^2).$$

556 Since $\sum_{i \in [2n]} \mu'_i$ is a constant, minimizing $\mathcal{L}_{\text{rec}}(\mathbf{W}; \mathbf{\Phi}, \mathbf{\Psi})$ is equivalent to maximizing

$$\mathcal{J}(\mathbf{W}) := \sum_{i \in [2n]} \mu'_i (\theta_i^\top \hat{\mathbf{W}}_i)^2 = \sum_{i \in [2n]} \hat{\mathbf{W}}_i^\top \mathbf{M}_i \hat{\mathbf{W}}_i.$$

557 For each $j \in [m]$, define $\mathbb{A}_j := \{i \in [2n] : \hat{\mathbf{W}}_i = \mathbf{W}_j\}$. Then $(\mathbb{A}_1, \dots, \mathbb{A}_m)$ is a partition of $[2n]$,
 558 and

$$\mathcal{J}(\mathbf{W}) = \sum_{j \in [m]} \mathbf{W}_j^\top \left(\sum_{i \in \mathbb{A}_j} \mathbf{M}_i \right) \mathbf{W}_j.$$

559 For a fixed partition, the objective decouples over j , and for each $j \in [m]$,

$$\max_{\|\mathbf{W}_j\|_2=1} \mathbf{W}_j^\top \left(\sum_{i \in \mathbb{A}_j} \mathbf{M}_i \right) \mathbf{W}_j = \lambda_{\max} \left(\sum_{i \in \mathbb{A}_j} \mathbf{M}_i \right),$$

560 attained by the top eigenvector. Since every choice of $(\mathbf{W}_1, \dots, \mathbf{W}_m)$ induces a partition of $[2n]$,
 561 the global minimization of the loss is equivalent to maximizing $\sum_{j \in [m]} \lambda_{\max} \left(\sum_{i \in \mathbb{A}_j} \mathbf{M}_i \right)$ over all
 562 partitions of $[2n]$ into m subsets. For the maximizing partition $(\mathbb{A}_1, \dots, \mathbb{A}_m)$, taking $\hat{\mathbf{W}}_j$ to be a
 563 unit-norm top eigenvector of $\sum_{i \in \mathbb{A}_j} \mathbf{M}_i$ yields a global minimizer $\hat{\mathbf{W}}$. Substituting back into (8),
 564 together with $\sum_{i \in [2n]} \mu'_i = 2 \sum_{i \in [n]} \mu'_i$, gives the stated reconstruction error:

$$\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \Phi) + \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \Psi) = s^{n-1}(1-s) \left(2 \sum_{i \in [n]} \mu'_i - \sum_{j \in [m]} \lambda_{\max} \left(\sum_{i \in \mathbb{A}_j} \mathbf{M}_i \right) \right) + o(1-s).$$

565 □

566 **Corollary 1.** *Suppose an SAE has $m \geq 2n$ latent coordinates and is trained by minimizing*
 567 $\mathcal{L}_{\text{rec}}(\mathbf{W}; \Phi) + \mathcal{L}_{\text{rec}}(\mathbf{W}; \Psi)$. *As the sparse factor s in (2) approaches 1, for any permutation*
 568 *matrix $\mathbf{P} \in \mathbb{R}^{m \times m}$, the weight matrix $\hat{\mathbf{W}} := [\Phi \ \Psi \ \mathbf{0}_{d \times (m-2n)}] \mathbf{P}$ is a global minimizer, and*
 569 $\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \Phi) + \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \Psi) = 0$.

570 *Proof.* Although Theorem 1 is stated for the regime $m < 2n$, the same argument extends beyond
 571 this restriction. When $m \geq 2n$, it gives the present corollary. We give a direct proof to make the
 572 construction explicit.

573 For $\mathbf{z} \in \mathbb{R}_+^n$, define the event $\mathbb{S}_k(\mathbf{z}) := \{\#\{i \in [n] : z_i = 0\} = k\}$ for $k \in \{0\} \cup [n]$, which denotes
 574 that exactly k entries of \mathbf{z} are zero. From (2), $\Pr(\mathbb{S}_k(\mathbf{z})) = \binom{n}{k} s^k (1-s)^{n-k}$. In the sparsity regime
 575 $s \rightarrow 1$, the loss admits the expansion

$$\mathcal{L}_{\text{rec}}(\mathbf{W}; \Phi) = s^{n-1}(1-s) \sum_{i \in [n]} \mu'_i \|\phi_i - \mathbf{W} \sigma(\mathbf{W}^\top \phi_i)\|_2^2 + o(1-s),$$

576 where $\mu'_i := \mathbb{E}[z_i^2 \mid z_i \neq 0] > 0$ for $i \in [n]$. An analogous expansion holds for $\mathcal{L}_{\text{rec}}(\mathbf{W}; \Psi)$. Adding
 577 the two and defining $\theta_i := \phi_i$, $\theta_{n+i} := \psi_i$ for $i \in [n]$, with $\mu'_{n+i} := \mu'_i$, we obtain

$$\mathcal{L}_{\text{rec}}(\mathbf{W}; \Phi, \Psi) = s^{n-1}(1-s) \sum_{i \in [2n]} \mu'_i \|\theta_i - \mathbf{W} \sigma(\mathbf{W}^\top \theta_i)\|_2^2 + o(1-s). \quad (9)$$

578 Since $s^{n-1}(1-s) > 0$ and $\mu'_i > 0$, minimizing $\mathcal{L}_{\text{rec}}(\mathbf{W}; \Phi, \Psi)$ in the regime $s \rightarrow 1$ is asymptotically
 579 equivalent to enforcing

$$\theta_i = \mathbf{W} \sigma(\mathbf{W}^\top \theta_i) \quad \text{for all } i \in [2n]. \quad (10)$$

580 For any permutation matrix $\mathbf{P} \in \mathbb{R}^{m \times m}$, substituting $\mathbf{W} = [\Phi \ \Psi \ \mathbf{0}_{d \times (m-2n)}] \mathbf{P}$ into (10) gives,
 581 for all $i \in [2n]$,

$$[\Phi \ \Psi \ \mathbf{0}_{d \times (m-2n)}] \mathbf{P} \sigma(\mathbf{P}^\top [\Phi \ \Psi \ \mathbf{0}_{d \times (m-2n)}]^\top \theta_i) = [\Phi \ \Psi] \sigma([\Phi \ \Psi]^\top \theta_i) = \theta_i,$$

582 which shows that $\hat{\mathbf{W}}$ is a global minimizer with $\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \Phi, \Psi) = 0$. □

583 **Proposition 3** (Complete statement of Proposition 1). *Consider the case $n = 1$, where the feature*
584 *directions $\phi, \psi \in \mathbb{R}^d$ satisfy $\rho := \phi^\top \psi \in (0, 1)$. Suppose the SAE has $m \geq 2$ latent coordinates,*
585 *and each column of \mathbf{W} has unit norm. Consider the loss [25]*

$$\mathcal{L}_{\text{rec}}(\mathbf{W}; \phi) + \mathcal{L}_{\text{rec}}(\mathbf{W}; \psi) + \lambda \mathbb{E}_z \left[\sum_{j \in [m]} \sqrt{[\sigma(\mathbf{W}^\top \phi z)]_j^2 + [\sigma(\mathbf{W}^\top \psi z)]_j^2} \right].$$

586 Define $\lambda^\dagger := \frac{\mathbb{E}[z^2]}{\mathbb{E}[z]}$, $\lambda^*(\rho) := \frac{(1-\rho)\lambda^\dagger}{2-\sqrt{1+\rho}}$, and $\lambda^\ddagger(\rho) := \sqrt{1+\rho} \lambda^\dagger$.

587 Up to any permutation matrix \mathbf{P} , a global minimizer $\hat{\mathbf{W}}$ of the loss is:

- 588 • If $\lambda < \lambda^*(\rho)$, then $\hat{\mathbf{W}} = [\phi \ \psi \ \mathbf{0}_{d \times (m-2)}] \mathbf{P}$, where $\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \phi) = \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \psi) = 0$.
- 589 • If $\lambda \in (\lambda^*(\rho), \lambda^\ddagger(\rho))$, then $\hat{\mathbf{W}} = \left[\frac{\phi + \psi}{\|\phi + \psi\|} \ \mathbf{0}_{d \times (m-1)} \right] \mathbf{P}$, where $\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \phi) = \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \psi) =$
590 $\frac{1-\rho}{2} \mathbb{E}[z^2]$.
- 591 • If $\lambda > \lambda^\ddagger(\rho)$, then $\hat{\mathbf{W}}$ satisfies $\sigma(\mathbf{W}^\top \phi) = \sigma(\mathbf{W}^\top \psi) = \mathbf{0}$, where $\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \phi) = \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \psi) =$
592 $\mathbb{E}[z^2]$.

593 *Proof.* Since $z \geq 0$, we have $\sigma(\mathbf{W}^\top \phi z) = z \sigma(\mathbf{W}^\top \phi)$ and $\sigma(\mathbf{W}^\top \psi z) = z \sigma(\mathbf{W}^\top \psi)$. Hence

$$\begin{aligned} \mathcal{L}_{\text{rec}}(\mathbf{W}; \phi) + \mathcal{L}_{\text{rec}}(\mathbf{W}; \psi) &= \mathbb{E}[z^2] \left(\|\phi - \mathbf{W} \sigma(\mathbf{W}^\top \phi)\|^2 + \|\psi - \mathbf{W} \sigma(\mathbf{W}^\top \psi)\|^2 \right), \\ \mathbb{E}_z \left[\sum_{j \in [m]} \sqrt{[\sigma(\mathbf{W}^\top \phi z)]_j^2 + [\sigma(\mathbf{W}^\top \psi z)]_j^2} \right] &= \mathbb{E}[z] \sum_{j \in [m]} \sqrt{[\sigma(\mathbf{W}^\top \phi)]_j^2 + [\sigma(\mathbf{W}^\top \psi)]_j^2}. \end{aligned}$$

594 Setting $\mathbf{u} := \sigma(\mathbf{W}^\top \phi)$, $\mathbf{v} := \sigma(\mathbf{W}^\top \psi)$, and $\tilde{\lambda} := \lambda \mathbb{E}[z]/\mathbb{E}[z^2]$, dividing through by $\mathbb{E}[z^2]$ reduces
595 the problem to minimizing

$$\mathcal{J}(\mathbf{W}) := \|\phi - \mathbf{W} \mathbf{u}\|^2 + \|\psi - \mathbf{W} \mathbf{v}\|^2 + \tilde{\lambda} \sum_{j \in [m]} \sqrt{[\mathbf{u}]_j^2 + [\mathbf{v}]_j^2}.$$

596 If ϕ and ψ activate different columns $j_a \neq j_b$, setting $[\mathbf{W}]_{[:,j_a]} = \phi$ and $[\mathbf{W}]_{[:,j_b]} = \psi$ gives zero
597 reconstruction error and group penalty 2; hence

$$L_{\text{sep}} := \mathcal{J}(\mathbf{W}) = 2\tilde{\lambda}.$$

598 If instead both activate the same unit column \mathbf{w} , writing $a := \mathbf{w}^\top \phi \geq 0$ and $b := \mathbf{w}^\top \psi \geq 0$,

$$\mathcal{J}([\mathbf{w}]) = 2 - (a^2 + b^2) + \tilde{\lambda} \sqrt{a^2 + b^2}.$$

599 With $\mathbf{M} := \phi \phi^\top + \psi \psi^\top$, we have $a^2 + b^2 = \mathbf{w}^\top \mathbf{M} \mathbf{w}$. Under $\rho \in (0, 1)$ the top eigenvector of \mathbf{M}
600 is

$$\hat{\mathbf{w}} := \frac{\phi + \psi}{\|\phi + \psi\|}, \quad a = b = \sqrt{\frac{1+\rho}{2}},$$

601 with eigenvalue $1 + \rho$. By the Rayleigh quotient, the feasible range of $a^2 + b^2$ over unit \mathbf{w} satisfies
602 $a^2 + b^2 \leq 1 + \rho$, with equality at $\hat{\mathbf{w}}$. Since $\mathcal{J}([\mathbf{w}])$ depends on \mathbf{w} through $r := \sqrt{a^2 + b^2}$, define

$$f(r) := 2 - r^2 + \tilde{\lambda} r, \quad r \in [0, \sqrt{1+\rho}].$$

603 Since f is concave in r , its minimum over $[0, \sqrt{1+\rho}]$ is attained at an endpoint. Evaluating at $r = 0$
604 and $r = \sqrt{1+\rho}$ gives

$$f(0) = 2, \quad f(\sqrt{1+\rho}) = (1-\rho) + \tilde{\lambda} \sqrt{1+\rho}.$$

605 Thus

$$L_{\text{sh}} := \min_{\mathbf{w}} \mathcal{J}([\mathbf{w}]) = \min \{2, (1-\rho) + \tilde{\lambda} \sqrt{1+\rho}\}.$$

606 Finally, consider the inactive solution where $\sigma(\mathbf{W}^\top \phi) = \mathbf{0}$ and $\sigma(\mathbf{W}^\top \psi) = \mathbf{0}$. In this case,

$$L_{\text{dead}} := \mathcal{J}(\mathbf{W}) = \|\phi\|^2 + \|\psi\|^2 = 2.$$

607 We now compare the three values:

$$L_{\text{sep}} = 2\tilde{\lambda}, \quad L_{\text{sh}} = \min \{2, (1 - \rho) + \tilde{\lambda}\sqrt{1 + \rho}\}, \quad L_{\text{dead}} = 2.$$

608 If $\tilde{\lambda} < \frac{1-\rho}{2-\sqrt{1+\rho}}$, then $2\tilde{\lambda} < (1 - \rho) + \tilde{\lambda}\sqrt{1 + \rho}$ and $2\tilde{\lambda} < 2$; hence $L_{\text{sep}} < L_{\text{sh}}$ and $L_{\text{sep}} < L_{\text{dead}}$.

609 If $\frac{1-\rho}{2-\sqrt{1+\rho}} < \tilde{\lambda} < \sqrt{1 + \rho}$, then $(1 - \rho) + \tilde{\lambda}\sqrt{1 + \rho} < 2\tilde{\lambda}$ and $(1 - \rho) + \tilde{\lambda}\sqrt{1 + \rho} < 2$; hence

610 $L_{\text{sh}} < L_{\text{sep}}$ and $L_{\text{sh}} < L_{\text{dead}}$. If $\tilde{\lambda} > \sqrt{1 + \rho}$, then $2 < 2\tilde{\lambda}$ and $2 < (1 - \rho) + \tilde{\lambda}\sqrt{1 + \rho}$; hence

611 $L_{\text{dead}} < L_{\text{sep}}$ and $L_{\text{dead}} < L_{\text{sh}}$.

612 Substituting $\tilde{\lambda} = \lambda \mathbb{E}[z]/\mathbb{E}[z^2]$ yields the result. When the shared solution is optimal, the reconstruction errors are

$$\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \phi) = \mathbb{E}[z^2](1 - a^2) = \frac{1-\rho}{2} \mathbb{E}[z^2], \quad \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \psi) = \mathbb{E}[z^2](1 - b^2) = \frac{1-\rho}{2} \mathbb{E}[z^2],$$

614 and when the inactive solution is optimal,

$$\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \phi) = \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \psi) = \mathbb{E}[z^2].$$

615

□

616 **Proposition 4.** Consider the case $n = 1$, where the feature directions $\phi, \psi \in \mathbb{R}^d$ satisfy $\rho :=$
 617 $\phi^\top \psi \in (0, 1)$. Suppose the SAE has $m \geq 2$ latent coordinates, and each column of \mathbf{W} has unit
 618 norm. Consider the loss [14]

$$\mathcal{L}_{\text{rec}}(\mathbf{W}; \phi) + \mathcal{L}_{\text{rec}}(\mathbf{W}; \psi) - \lambda \mathbb{E}_z \left[\sigma(\mathbf{W}^\top \phi z)^\top \sigma(\mathbf{W}^\top \psi z) \right].$$

619 Define $\lambda^*(\rho) := \frac{2(1-\rho)}{1+\rho}$.

620 Up to any permutation matrix \mathbf{P} , a global minimizer $\hat{\mathbf{W}}$ of the loss is:

621 • If $\lambda < \lambda^*(\rho)$, then $\hat{\mathbf{W}} = [\phi \ \psi \ \mathbf{0}_{d \times (m-2)}] \mathbf{P}$, where $\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \phi) = \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \psi) = 0$.

622 • If $\lambda > \lambda^*(\rho)$, then $\hat{\mathbf{W}} = \left[\frac{\phi + \psi}{\|\phi + \psi\|} \ \mathbf{0}_{d \times (m-1)} \right] \mathbf{P}$, where $\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \phi) = \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \psi) = \frac{1-\rho}{2} \mathbb{E}[z^2]$.

623 *Proof.* Since $z \geq 0$, we have $\sigma(\mathbf{W}^\top \phi z) = z \sigma(\mathbf{W}^\top \phi)$ and $\sigma(\mathbf{W}^\top \psi z) = z \sigma(\mathbf{W}^\top \psi)$. Hence

$$\begin{aligned} \mathcal{L}_{\text{rec}}(\mathbf{W}; \phi) + \mathcal{L}_{\text{rec}}(\mathbf{W}; \psi) &= \mathbb{E}[z^2] \left(\|\phi - \mathbf{W} \sigma(\mathbf{W}^\top \phi)\|^2 + \|\psi - \mathbf{W} \sigma(\mathbf{W}^\top \psi)\|^2 \right), \\ \mathbb{E}_z \left[\sigma(\mathbf{W}^\top \phi z)^\top \sigma(\mathbf{W}^\top \psi z) \right] &= \mathbb{E}[z^2] \sigma(\mathbf{W}^\top \phi)^\top \sigma(\mathbf{W}^\top \psi). \end{aligned}$$

624 Setting $\mathbf{u} := \sigma(\mathbf{W}^\top \phi)$ and $\mathbf{v} := \sigma(\mathbf{W}^\top \psi)$, the loss factors as $\mathbb{E}[z^2] \mathcal{J}(\mathbf{W})$, where

$$\mathcal{J}(\mathbf{W}) := \|\phi - \mathbf{W} \mathbf{u}\|^2 + \|\psi - \mathbf{W} \mathbf{v}\|^2 - \lambda \mathbf{u}^\top \mathbf{v}.$$

625 If ϕ and ψ activate different columns $j_a \neq j_b$, setting $[\mathbf{W}]_{[:,j_a]} = \phi$ and $[\mathbf{W}]_{[:,j_b]} = \psi$ gives zero
 626 reconstruction error and $\mathbf{u}^\top \mathbf{v} = 0$; hence

$$L_{\text{sep}} := \mathcal{J}(\mathbf{W}) = 0.$$

627 If instead both activate the same unit column \mathbf{w} , writing $a := \mathbf{w}^\top \phi \geq 0$ and $b := \mathbf{w}^\top \psi \geq 0$,

$$\mathcal{J}([\mathbf{w}]) = 2 - a^2 - b^2 - \lambda ab = 2 - \mathbf{w}^\top \mathbf{M}_\lambda \mathbf{w},$$

628 where $\mathbf{M}_\lambda := \phi \phi^\top + \psi \psi^\top + \frac{\lambda}{2} (\phi \psi^\top + \psi \phi^\top)$. Under $\rho \in (0, 1)$ and $\lambda \geq 0$, the top eigenvector is

$$\hat{\mathbf{w}} := \frac{\phi + \psi}{\|\phi + \psi\|}, \quad a = b = \sqrt{\frac{1+\rho}{2}}.$$

629 By the Rayleigh quotient, $\mathbf{w}^\top \mathbf{M}_\lambda \mathbf{w} \leq (1 + \rho)(1 + \frac{\lambda}{2})$ for unit \mathbf{w} , with equality at $\hat{\mathbf{w}}$. Therefore

$$L_{\text{sh}} := \min_{\mathbf{w}} \mathcal{J}([\mathbf{w}]) = 2 - (1 + \rho)(1 + \frac{\lambda}{2}) = (1 - \rho) - \lambda \frac{1 + \rho}{2}.$$

630 Finally, consider the inactive solution where $\sigma(\mathbf{W}^\top \phi) = \mathbf{0}$ and $\sigma(\mathbf{W}^\top \psi) = \mathbf{0}$. In this case,

$$L_{\text{dead}} := \mathcal{J}(\mathbf{W}) = \|\phi\|^2 + \|\psi\|^2 = 2.$$

631 We now compare the three values:

$$L_{\text{sep}} = 0, \quad L_{\text{sh}} = (1 - \rho) - \lambda \frac{1 + \rho}{2}, \quad L_{\text{dead}} = 2.$$

632 Since $\lambda \geq 0$, we have $L_{\text{dead}} = 2 > 0 = L_{\text{sep}}$, so the dead solution is never optimal. Comparing the
633 remaining two, $L_{\text{sh}} < L_{\text{sep}}$ if and only if $\lambda > \frac{2(1 - \rho)}{1 + \rho} = \lambda^*(\rho)$.

634 When $\lambda > \lambda^*(\rho)$, the reconstruction errors at the shared minimizer are

$$\mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \phi) = \mathbb{E}[z^2](1 - a^2) = \frac{1 - \rho}{2} \mathbb{E}[z^2], \quad \mathcal{L}_{\text{rec}}(\hat{\mathbf{W}}; \psi) = \mathbb{E}[z^2](1 - b^2) = \frac{1 - \rho}{2} \mathbb{E}[z^2].$$

635 □

636 E Experimental Details

637 This section provides experimental details for the empirical analyses in Sections 3 and 6.

638 All experiments are conducted on a single NVIDIA A100 GPU.

639 E.1 Experimental Details for Measuring Cross-Modal Feature Heterogeneity

640 This section provides experimental details for Figure 2 in Section 3, which reports the distribution
641 of cosine distances between estimated image and text feature directions grouped by coactivation
642 correlation. Recall that we use decoder columns $\hat{\phi}_i$ and $\hat{\psi}_j$ of modality-specific SAEs as estimates
643 of image and text feature directions, and that the coactivation correlation $c_{i,j} := [\mathbf{C}]_{i,j}$ between the
644 i -th image latent and the j -th text latent on paired embeddings (5) serves as a proxy for semantic
645 correspondence. For each pair (i, j) , we measure the cosine distance between $\hat{\phi}_i$ and $\hat{\psi}_j$.

646 We group all m^2 index pairs $(i, j)^3$ by their correlation value $c_{i,j}$ and examine the distribution
647 of cosine distances between the corresponding feature vectors. A larger $c_{i,j}$ indicates that the
648 corresponding image and text features are more likely to represent the same shared concept. Under
649 perfect cross-modal feature alignment, their cosine distances would concentrate near zero. We train
650 modality-specific SAEs with latent dimension $m = 8192$, using the Top- K [33] activation function
651 with $K = 8$. We follow the training protocol of Papadimitriou et al. [38]. Specifically, we use the
652 AdamW optimizer with learning rate 5×10^{-4} and weight decay 10^{-5} , batch size 1024, and a cosine
653 schedule with 5% linear warmup over 30 training epochs.

654 E.2 Implementation Details for Synthetic Data Experiments

655 This section provides the data-generation specification, SAE architectures, and optimization hyperpa-
656 rameters used in Section 6.1.

657 **Dataset.** Following prior work [19, 49], we generate each pair of synthetic embeddings in two steps.
658 We first sample a sparse latent code \mathbf{z} from a Bernoulli–Exponential process, and then form (\mathbf{x}, \mathbf{y})
659 from the ground-truth feature matrices (Φ, Ψ) via (1). The embedding dimension is $d = 256$, and
660 the feature matrices $\Phi, \Psi \in \mathbb{R}^{d \times n}$ contain $n = n_S + n_I + n_T = 2048$ columns in total: $n_S = 1024$
661 shared concepts present in both modalities, $n_I = 512$ image-only concepts (with $\psi_i = \mathbf{0}$), and
662 $n_T = 512$ text-only concepts (with $\phi_i = \mathbf{0}$). For each shared concept $i \in [n_S]$, we sample ϕ_i and
663 ψ_i as unit-norm vectors satisfying $\cos(\phi_i, \psi_i) = \alpha$ for a prescribed cross-modal feature alignment
664 α , giving us direct control over the latent structure of the generative process.

³We exclude dead latent coordinates, which are zero across all inputs, so fewer than m^2 pairs remain in practice.

665 Modality-specific feature directions are sampled as unit-norm vectors with maximum pairwise
666 interference bounded by $\epsilon_{\max} = 0.30$. Each coordinate of the latent code $\mathbf{z} \in \mathbb{R}_+^n$ is independently
667 set to zero with probability $s = 0.99$, and otherwise drawn i.i.d. from $\text{Exp}(\beta)$ with rate $\beta = 1$. The
668 paired embeddings are then formed as $\mathbf{x} = \Phi\mathbf{z} + \epsilon_I$ and $\mathbf{y} = \Psi\mathbf{z} + \epsilon_T$ via (1), with independent
669 observation noise $\epsilon_I, \epsilon_T \sim \mathcal{N}(\mathbf{0}, \sigma_{\text{obs}}^2 \mathbf{I})$ of standard deviation $\sigma_{\text{obs}} = 0.05$. We generate 50,000
670 paired embeddings for training and 10,000 for evaluation, and repeat each experiment over three
671 independent runs, reporting the average.

672 **Training.** The shared SAE has latent dimension $m = 8192$ and uses the Top- K activation function
673 [33] with $K = 16$. Each modality-specific SAE uses latent dimension $m = 4096$, so that the
674 total number of learnable parameters matches that of the shared SAE. The two baselines [14, 25]
675 augment the shared SAE by training with an auxiliary loss.

676 **Optimization.** All SAEs are trained with AdamW ($\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$, weight decay
677 0) at a constant learning rate of 5×10^{-4} with no warmup. We use a batch size of 256 and train for
678 10 epochs, corresponding to approximately 1,950 optimization steps.

679 E.3 Evaluation Metrics for Synthetic Data Experiments

680 We provide formal definitions of the five metrics introduced in Section 6.1.

681 We assess each method using these metrics: (i) and (ii) measure how well the SAE reconstructs
682 embeddings and recovers features, (iii) and (iv) measure how well it aligns latent codes across
683 modalities, and (v) measures whether cross-modal features collapse into a single feature.

- 684 (i) *Reconstruction Error.* We measure the mean squared error (MSE) between input embeddings
685 \mathbf{x} (or \mathbf{y}) and their reconstructions $\tilde{\mathbf{x}}$ (or $\tilde{\mathbf{y}}$) on an evaluation set, averaged over both modalities.
- 686 (ii) *Feature Recovery Error.* Similar to Reconstruction Error, but using the feature direction vectors
687 ϕ_i (or ψ_i) as inputs, we measure their reconstruction MSE averaged over both modalities.
- 688 (iii) *Alignment Error.* We measure the cosine distance between the image and text latent codes
689 $\tilde{\mathbf{z}}_I(\mathbf{x})$ and $\tilde{\mathbf{z}}_T(\mathbf{y})$ obtained from paired embeddings (\mathbf{x}, \mathbf{y}) on an evaluation set.
- 690 (iv) *Feature Alignment Error.* Similar to Alignment Error, but using features of the same concept
691 (ϕ_i, ψ_i) as inputs, we measure the cosine distance between the latent codes $\tilde{\mathbf{z}}_I(\phi_i)$ and $\tilde{\mathbf{z}}_T(\psi_i)$.
- 692 (v) *(Cross-Modal) Feature Collapse Rate.* We measure the fraction of feature pairs (ϕ_i, ψ_i) that
693 represent the same concept across modalities and are assigned to the same latent coordinate,
694 indicating whether the SAE merges two modality-specific features into a single learned feature.

695 Note that metrics (ii), (iv), and (v) require access to the ground-truth features (Φ, Ψ) and thus apply
696 only in the synthetic setting. The mathematical definitions are given below.

697 (i) *Reconstruction Error:*

$$\frac{1}{2} \mathbb{E} \left[\|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2 + \|\mathbf{y} - \tilde{\mathbf{y}}\|_2^2 \right] = \frac{1}{2} \mathbb{E} \left[\|\mathbf{x} - \mathbf{W}\sigma(\mathbf{W}^\top \mathbf{x})\|_2^2 + \|\mathbf{y} - \mathbf{W}\sigma(\mathbf{W}^\top \mathbf{y})\|_2^2 \right].$$

698 (ii) *Feature Recovery Error:*

$$\frac{1}{2n_S} \sum_{i \in [n_S]} \left(\|\psi_i - \mathbf{W}\sigma(\mathbf{W}^\top \psi_i)\|_2^2 + \|\phi_i - \mathbf{W}\sigma(\mathbf{W}^\top \phi_i)\|_2^2 \right).$$

699 (iii) *Alignment Error:*

$$1 - \mathbb{E} [\cos(\tilde{\mathbf{z}}(\mathbf{x}), \tilde{\mathbf{z}}(\mathbf{y}))] = 1 - \mathbb{E} [\cos(\sigma(\mathbf{W}^\top \mathbf{x}), \sigma(\mathbf{W}^\top \mathbf{y}))].$$

700 (iv) *Feature Alignment Error:*

$$1 - \frac{1}{n_S} \sum_{i \in [n_S]} \cos(\sigma(\mathbf{W}^\top \phi_i), \sigma(\mathbf{W}^\top \psi_i)).$$

701 (v) *(Cross-Modal) Feature Collapse Rate:*

$$\frac{1}{n_S} \sum_{i \in [n_S]} \mathbf{1} \left[\arg \max_{j \in [m]} \cos([\mathbf{W}]_{:,j}, \phi_i) = \arg \max_{j \in [m]} \cos([\mathbf{W}]_{:,j}, \psi_i) \right] \in [0, 1].$$

702 **E.4 Implementation Details for Real-World Experiments**

703 **Implementation Details.** We use the CC-3M [48] dataset with approximately 2.82M image-caption
 704 pairs after dropping invalid URLs. We optimize all SAEs with AdamW ($\beta_1 = 0.9, \beta_2 = 0.999$)
 705 at learning rate 5×10^{-4} , weight decay 10^{-5} , batch size 1024, gradient clipping at norm 1.0,
 706 and a cosine schedule with 5% linear warmup followed by cosine decay to zero. We use latent
 707 dimension $m = 8192$ and Top- K activation [33] with $K = 32$, and we adopt the training protocol
 708 of Papadimitriou et al. [38].

709 **Licenses of datasets.** We use CC-3M [48] and MS-COCO [29] for our experiments. CC-3M is released
 710 by Google for free use under the terms of its repository license. MS-COCO annotations are released
 711 under CC BY 4.0.

712 **F Additional Experimental Results**

713 **F.1 Evidence for Cross-Modal Feature Heterogeneity on Larger Backbones**

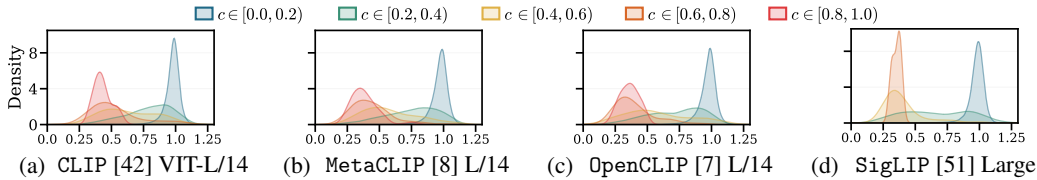


Figure 9: Distribution of cosine distances between image-text feature pairs estimated from embeddings of four VLMs. We group image-text feature pairs by their co-activation correlation $c_{i,j}$ in (5), shown in different colors. The value $c_{i,j}$ measures how strongly the i -th image feature and the j -th text feature co-activate on paired image-text embeddings, so pairs with larger correlations are more likely to represent the same shared concept. Across all models, the distribution does not concentrate near 0 but remains centered around a positive value, even for high-correlation pairs ($c \geq 0.8$, shown in blue). This observation supports the presence of cross-modal feature heterogeneity.

714 Figure 2 in the main text uses the base variants (ViT-B/32) of the four VLMs. To check whether
 715 the observed cross-modal feature heterogeneity persists at larger scales, we repeat the same analysis
 716 on the corresponding large variants (ViT-L/14) of CLIP, MetaCLIP, OpenCLIP, and SigLIP2, and
 717 report the results in Figure 9. The overall trend is consistent with Figure 2: feature pairs with larger
 718 coactivation correlation tend to have smaller cosine distances, but they do not concentrate near zero.
 719 This indicates that cross-modal feature heterogeneity is not an artifact of model capacity and persists
 720 across backbone sizes.

721 **F.2 Evidence of Feature Collapse on Real-World Data**

722 We complement the synthetic results in Figure 5a with a real-world experiment that checks whether
 723 the feature collapse predicted by Theorem 1 actually occurs in SAEs trained on VLM embeddings.

724 In principle, SAEs use a large latent dimension m , so one might expect to be in the safe regime
 725 of Corollary 1 ($m \geq 2n$). Practically, however, a substantial fraction of latent coordinates are
 726 dead [50, 18], so the effective number of usable coordinates can be much smaller than m . This
 727 pushes the SAE toward the regime of Theorem 1, where directionally close cross-modal features are
 728 collapsed into a single feature. We conduct a simple diagnostic to detect such collapses on real data.

729 **Setup.** We follow the same protocol as Appendix E.1 on CLIP ViT-B/32 embeddings of
 730 MS-COCO [29], except that we now train a single *shared* SAE on the union of image and text
 731 embeddings, with $m = 8192$ and Top- K activation [33] ($K = 8$).

732 **Diagnostic.** We compute the coactivation correlation matrix C in (5) from the shared encoder, and
 733 obtain the Hungarian matching \hat{P} in (6). The two regimes in Section 4 make opposite predictions
 734 about \hat{P} :

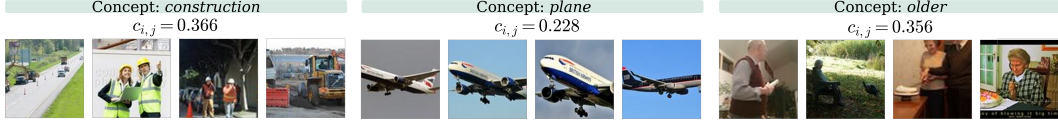


Figure 10: Examples of images that activates collapsed image-text feature pairs identified by a single shared SAE on CLIP embeddings of MS-COCO [29]. We define a coordinate is *collapsed* when $\hat{\mathbf{P}}(i) = i$ in the Hungarian matching of (6) applied to the correlation matrix \mathbf{C} in (5), indicating that the same latent coordinate represents a concept in both modalities. Each row shows the top-activating images for one collapsed coordinate, along with its coactivation correlation $c_{i,j}$. We show the three collapsed pairs with the *lowest* $c_{i,j}$ rather than the highest. Even in this worst case, the coordinates correspond to coherent concepts (*plane*, *older*, *construction*), but $c_{i,j}$ remains low (0.23–0.37).

- 735 • Under modality split (Corollary 1), the same concept activates two *different* coordinates across
 736 modalities, so a Hungarian match aligns the coordinate with index i to a different coordinate
 737 $\hat{\mathbf{P}}(i) \neq i$.
- 738 • Under collapse (Theorem 1), the same concept activates the *same* coordinate across modalities, so
 739 the match is to itself, $\hat{\mathbf{P}}(i) = i$.
- 740 This leads to a natural diagnostic for collapse:

$$\text{collapse rate} := \frac{|\{i \in [\tilde{m}] : \hat{\mathbf{P}}(i) = i\}|}{\tilde{m}},$$

741 where \tilde{m} is the number of live coordinates after excluding dead neurons.

742 **Result.** On the embeddings above, we observe a collapse rate of 4.64%. While small in absolute
 743 terms, this confirms that the collapse regime characterized by Theorem 1 does occur in real-world data.
 744 To verify that the diagonal matches reflect genuine semantic correspondences rather than coincidences,
 745 Figure 10 visualizes top-activating images for three collapsed coordinates. We deliberately select the
 746 *three collapsed pairs with the lowest* $c_{i,j}$, presenting a worst case rather than a favorable one. Even
 747 in this worst case, each coordinate corresponds to a coherent concept (*plane*, *older*, *construction*).

748 A further observation is that even at these collapsed coordinates, the coactivation correlation $c_{i,j}$ stays
 749 modest (0.23–0.37). This may seem counterintuitive, since collapsed pairs share a single column and
 750 one might expect strong coactivation. But it is consistent with the analysis in Proposition 1. When
 751 image and text features for the same concept are directionally distinct yet forced to share one column,
 752 the optimal shared direction is the compromise $\frac{\phi_i + \psi_i}{\|\phi_i + \psi_i\|}$ (see Propositions 4 and 3 for details), which
 753 fits neither modality well and produces only weak coactivation. In other words, collapse hurts not
 754 just reconstruction but cross-modal alignment as well.

755 F.3 Real-World Results Across Experimental Configurations

756 In the main experiment in Section 6.2, Table 1 reports results for a single configuration (CLIP
 757 ViT-B/32 embeddings with Top- K activation at $K = 32$). To verify that the conclusions drawn there
 758 are not artifacts of this particular setup, we conduct extended experiments along three orthogonal
 759 axes: (i) the sparsity level K in Top- K activation, (ii) the pre-trained VLM backbone (varying both
 760 architecture scale and pre-training corpus), and (iii) the choice of sparsifying activation function. All
 761 other settings follow Section 6.2, and we report averages over three independent runs. Tables 2, 3,
 762 and 4 report the per-configuration results. Across every configuration, our method matches the lowest
 763 reconstruction error of modality-specific SAEs while delivering the strongest cross-modal retrieval
 764 performance, mirroring the conclusions of Table 1.

765 **Robustness to sparsity level K .** We first vary the Top- K sparsity level on CLIP ViT-B/32 em-
 766 beddings (Table 2). Together with the $K = 32$ result in Table 1, this covers $K \in \{16, 32, 64\}$.
 767 Across all three sparsity levels, our method preserves the lowest reconstruction error inherited from
 768 modality-specific SAEs. On cross-modal retrieval, $K = 32$ stands out as the best-performing sparsity
 769 level: our method peaks there with I→T Recall@1 of 16.0, a level not matched by any method at
 770 either $K = 16$ or $K = 64$. The retrieval gain over the strongest baseline is most pronounced in the
 771 sparser regime: at $K = 16$, our method more than doubles the performance of the strongest baseline

Table 2: Effect of Top- K sparsity level K on CLIP ViT-B/32 embeddings ($L = 8192$). Each section corresponds to a different value of K , and the $K = 32$ setting is reported in Table 1. Best and second-best results within each section are shown in bold and underlined, respectively.

Methods	MS-COCO [29]						ImageNet1K [13]		
	Recon. (\downarrow)	Image-to-Text (\uparrow)			Text-to-Image (\uparrow)			Recon. (\downarrow)	Zero-shot (\uparrow)
	MSE	R@1	R@5	R@10	R@1	R@5	R@10	MSE	Accuracy
$K = 16$									
Shared SAE	0.131	3.0 (± 1.9)	7.7 (± 3.3)	10.9 (± 4.1)	3.0 (± 0.3)	7.8 (± 0.7)	11.3 (± 1.1)	0.169	12.6 (± 4.1)
+ Iso-Energy alignment loss	<u>0.131</u>	4.2 (± 1.2)	9.6 (± 2.8)	13.6 (± 3.8)	<u>3.1</u> (± 0.2)	8.2 (± 0.8)	12.1 (± 1.3)	<u>0.168</u>	15.3 (± 3.8)
+ Group-sparse loss	0.134	<u>5.1</u> (± 0.1)	<u>12.2</u> (± 0.2)	<u>17.2</u> (± 0.2)	<u>2.9</u> (± 0.1)	<u>8.8</u> (± 0.5)	<u>12.9</u> (± 0.5)	<u>0.170</u>	<u>20.3</u> (± 0.6)
Modality-Specific SAEs	0.127	0.0 (± 0.0)	0.1 (± 0.1)	0.2 (± 0.1)	0.0 (± 0.0)	0.1 (± 0.1)	0.2 (± 0.1)	0.165	0.1 (± 0.1)
+ Post-hoc Alignment (Ours)	0.127	13.3 (± 0.7)	29.4 (± 0.3)	38.9 (± 0.7)	9.6 (± 0.3)	23.0 (± 0.5)	31.8 (± 0.5)	0.165	22.9 (± 0.5)
$K = 64$									
Shared SAE	0.063	9.8 (± 1.3)	21.2 (± 2.0)	28.0 (± 2.9)	6.0 (± 0.6)	15.2 (± 1.2)	21.3 (± 1.5)	0.083	15.0 (± 1.5)
+ Iso-Energy alignment loss	0.063	9.0 (± 0.8)	19.9 (± 0.4)	26.7 (± 0.5)	5.3 (± 0.6)	13.4 (± 1.5)	18.8 (± 2.0)	0.084	16.3 (± 1.2)
+ Group-sparse loss	0.091	11.8 (± 0.5)	26.0 (± 0.8)	34.6 (± 1.1)	7.0 (± 0.3)	18.4 (± 0.5)	26.4 (± 0.5)	0.116	33.0 (± 0.4)
Modality-Specific SAEs	0.062	0.0 (± 0.0)	0.1 (± 0.0)	0.2 (± 0.1)	0.0 (± 0.0)	0.1 (± 0.0)	0.3 (± 0.1)	0.083	0.1 (± 0.0)
+ Post-hoc Alignment (Ours)	0.062	11.2 (± 0.8)	26.0 (± 1.6)	34.9 (± 2.2)	8.5 (± 0.8)	21.5 (± 2.1)	29.7 (± 2.8)	0.083	19.3 (± 1.2)

Table 3: Effect of pre-trained VLM backbone with $K = 32$ and Top- K activation. Each section corresponds to a different backbone, and the CLIP ViT-B/32 setting is reported in Table 1. Best and second-best results within each section are shown in bold and underlined, respectively.

Methods	MS-COCO [29]						ImageNet1K [13]		
	Recon. (\downarrow)	Image-to-Text (\uparrow)			Text-to-Image (\uparrow)			Recon. (\downarrow)	Zero-shot (\uparrow)
	MSE	R@1	R@5	R@10	R@1	R@5	R@10	MSE	Accuracy
<i>CLIP ViT-L/14</i> ($L = 12288$)									
Shared SAE	0.119	18.2 (± 0.2)	33.6 (± 0.4)	41.8 (± 0.7)	9.9 (± 0.8)	22.9 (± 0.9)	30.7 (± 1.2)	0.155	38.9 (± 0.6)
+ Iso-Energy alignment loss	0.119	16.8 (± 1.1)	31.2 (± 1.2)	39.3 (± 1.7)	9.4 (± 0.5)	21.3 (± 0.9)	28.9 (± 1.2)	<u>0.154</u>	<u>39.4</u> (± 0.8)
+ Group-sparse loss	0.137	14.2 (± 0.2)	29.6 (± 0.7)	38.5 (± 1.0)	8.2 (± 0.4)	20.3 (± 0.4)	28.4 (± 0.3)	0.176	45.0 (± 0.3)
Modality-Specific SAEs	0.117	0.0 (± 0.0)	0.1 (± 0.0)	0.2 (± 0.0)	0.0 (± 0.0)	0.1 (± 0.0)	0.2 (± 0.1)	0.154	0.1 (± 0.0)
+ Post-hoc Alignment (Ours)	0.117	22.6 (± 0.7)	43.6 (± 1.2)	54.9 (± 1.2)	15.3 (± 0.4)	33.0 (± 0.4)	42.9 (± 0.3)	0.154	35.9 (± 1.1)
<i>OpenCLIP ViT-B/32</i> ($L = 8192$)									
Shared SAE	0.116	8.9 (± 0.3)	18.5 (± 0.3)	24.4 (± 0.5)	4.8 (± 0.2)	12.2 (± 0.2)	17.4 (± 0.2)	0.149	18.7 (± 0.7)
+ Iso-Energy alignment loss	0.117	9.2 (± 0.8)	18.6 (± 1.3)	24.5 (± 1.2)	4.8 (± 0.5)	12.3 (± 0.8)	17.4 (± 1.1)	0.149	18.9 (± 1.3)
+ Group-sparse loss	0.141	<u>10.1</u> (± 0.4)	<u>23.1</u> (± 0.4)	<u>31.8</u> (± 0.5)	<u>6.1</u> (± 0.0)	<u>17.1</u> (± 0.2)	<u>25.4</u> (± 0.3)	0.176	34.3 (± 0.4)
Modality-Specific SAEs	0.115	0.0 (± 0.0)	0.1 (± 0.0)	0.2 (± 0.1)	0.0 (± 0.0)	0.1 (± 0.0)	0.2 (± 0.1)	0.149	0.1 (± 0.0)
+ Post-hoc Alignment (Ours)	0.115	21.0 (± 0.9)	41.2 (± 1.1)	52.0 (± 1.0)	11.2 (± 0.1)	26.5 (± 0.6)	36.1 (± 0.9)	0.149	29.2 (± 0.7)
<i>OpenCLIP ViT-L/14</i> ($L = 12288$)									
Shared SAE	0.133	26.5 (± 0.4)	45.5 (± 1.8)	54.5 (± 2.1)	15.8 (± 0.8)	32.0 (± 1.5)	40.6 (± 1.7)	0.164	50.9 (± 0.7)
+ Iso-Energy alignment loss	0.134	<u>27.1</u> (± 0.4)	<u>46.2</u> (± 1.1)	<u>55.5</u> (± 0.7)	<u>16.2</u> (± 0.5)	<u>32.9</u> (± 0.9)	<u>41.6</u> (± 0.8)	0.166	50.3 (± 0.5)
+ Group-sparse loss	0.162	23.7 (± 0.2)	43.3 (± 0.5)	53.9 (± 0.1)	13.1 (± 0.3)	29.1 (± 0.4)	38.6 (± 0.5)	0.196	53.9 (± 0.3)
Modality-Specific SAEs	0.129	0.1 (± 0.0)	0.1 (± 0.1)	0.2 (± 0.1)	0.0 (± 0.0)	0.1 (± 0.1)	0.2 (± 0.1)	0.164	0.1 (± 0.0)
+ Post-hoc Alignment (Ours)	0.129	27.4 (± 2.2)	50.6 (± 1.8)	62.0 (± 1.4)	18.1 (± 2.0)	37.6 (± 3.0)	47.9 (± 3.1)	0.164	42.4 (± 0.4)

772 on both retrieval directions. At $K = 64$, the auxiliary-loss baselines benefit from the larger active
773 latent budget and the group-sparse loss becomes competitive on I \rightarrow T Recall@1; however, retrieval
774 scores of our method also drop relative to $K = 32$, so the narrowed gap reflects a regime where
775 no method attains the peak achievable at $K = 32$, rather than one where the baselines genuinely
776 surpass us. Even at $K = 64$, our method still leads on T \rightarrow I Recall and for top-10 retrieval in both
777 directions, while maintaining a substantially lower reconstruction error. Overall, $K = 32$ is the
778 optimal operating point for cross-modal retrieval, and at this K our method generally outperforms all
779 baselines.

780 **Robustness across vision-language model backbones.** We next evaluate three additional back-
781 bones at $K = 32$ with Top- K activation, as shown in Table 3. We span two architecture scales
782 (ViT-B/32 vs. ViT-L/14) and two pre-training corpora (OpenAI CLIP [42] vs. OpenCLIP [7]). For
783 each backbone, we scale the SAE width m with the embedding dimension d so that the expansion ratio
784 m/d matches the setting of Table 1, ensuring a fair comparison across architectures of different sizes.
785 On all three backbones, our method matches the lowest reconstruction error and obtains the strongest
786 cross-modal retrieval performance, with the same qualitative behavior as the CLIP ViT-B/32 result
787 in Table 1. The retrieval gain over the strongest baseline is largest on the ViT-B/32 backbones and
788 smaller on ViT-L/14, reflecting that the L/14 embedding spaces already exhibit stronger cross-modal
789 alignment and thus leave less room for improvement at Recall@1. The advantage on T \rightarrow I Recall@1

Table 4: Effect of sparsifying activation function on CLIP ViT-B/32 ($L = 8192$, $K = 32$). The Top- K setting is reported in Table 1; this table reports the BatchTop- K [5] variant. Best and second-best results are shown in bold and underlined, respectively.

Methods	MS-COCO [29]						ImageNet1K [13]		
	Recon. (\downarrow)	Image-to-Text (\uparrow)			Text-to-Image (\uparrow)			Recon. (\downarrow)	Zero-shot (\uparrow)
	MSE	R@1	R@5	R@10	R@1	R@5	R@10	MSE	Accuracy
	<i>BatchTop-K</i>								
Shared SAE	0.089	6.3 (± 0.6)	15.9 (± 1.5)	22.0 (± 1.7)	3.9 (± 0.6)	11.3 (± 0.9)	16.8 (± 1.0)	<u>0.115</u>	16.4 (± 3.1)
+ Iso-Energy alignment loss	<u>0.089</u>	5.8 (± 0.8)	14.8 (± 0.8)	20.4 (± 0.9)	3.3 (± 1.4)	9.3 (± 3.3)	13.9 (± 3.9)	0.115	13.0 (± 1.8)
+ Group-sparse loss	0.105	<u>7.8</u> (± 0.3)	<u>19.0</u> (± 0.4)	<u>26.2</u> (± 0.4)	<u>4.9</u> (± 0.0)	<u>13.7</u> (± 0.1)	<u>20.6</u> (± 0.1)	0.132	28.0 (± 0.7)
Modality-Specific SAEs	0.088	0.0 (± 0.0)	0.0 (± 0.0)	0.1 (± 0.0)	0.0 (± 0.0)	0.1 (± 0.0)	0.2 (± 0.0)	0.114	0.1 (± 0.0)
+ Post-hoc Alignment (Ours)	0.088	16.2 (± 0.8)	34.6 (± 1.0)	44.7 (± 1.4)	10.5 (± 0.6)	25.5 (± 0.8)	35.0 (± 0.9)	0.114	<u>22.8</u> (± 1.0)

790 is nevertheless consistently positive across all three backbones. As in the main experiment, the only
791 metric where the group-sparse baseline occasionally surpasses our method is ImageNet zero-shot
792 accuracy, but this comes at the cost of substantially degraded reconstruction.

793 **Robustness to activation function.** Finally, we replace Top- K with BatchTop- K [5], which selects
794 the top- K activations across the entire batch rather than per token, on CLIP ViT-B/32 with $K = 32$,
795 as shown in Table 4. Our method again attains the lowest reconstruction error (0.0878) and the
796 strongest cross-modal retrieval, improving I \rightarrow T Recall@1 from 7.81 to 16.19 and T \rightarrow I Recall@1
797 from 4.88 to 10.46 over the strongest baseline. The magnitude of these gains closely matches the
798 corresponding Top- K setting in Table 1, confirming that the benefits of post-hoc alignment are not
799 specific to vanilla Top- K and transfer cleanly to alternative sparsifying activations.