
Feature Recovery Requires Structured Event Regimes in Sparse Reconstruction

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Sparse autoencoders are often used in mechanistic interpretability as if sparse
2 reconstruction should recover the features represented by a model. Recent work
3 shows that this recovery is fragile, but it remains unclear which failures come
4 from the SAE architecture, the encoder, optimization, or finite data. We show
5 that several failures can already be incentivized by the population-level sparse-
6 reconstruction objective. We study by how much residual mass projects above
7 the sparsity threshold in a positive linear latent-ray model; merging and splitting
8 arise as static properties of this objective, while absorption and seed-dependent
9 alternatives arise sequentially as earlier selections change the residual field. We
10 also separate two notions often conflated in interpretability practice: recovering a
11 ground-truth direction and recovering the activation pattern of that feature. Neither
12 implies the other in general; they coincide under specific structural conditions,
13 such as single-feature event dominance and regular simplex structure in learned
14 symmetric geometries. Sparse reconstruction therefore recovers ground-truth
15 features only in structured event regimes; outside them, the objective can favor non-
16 canonical but reconstruction-useful directions, and a one-ReLU encoder introduces
17 a further representability gap governed by whether the oracle gate is affine-ReLU
18 approximable. Overall, our results refine the existing analysis of SAE behavior and
19 provide a unified perspective on ground-truth feature recovery studies.

20 1 Introduction

21 Mechanistic interpretability often studies neural representations through the lens of features: direc-
22 tions or low-dimensional components that correspond to meaningful computational variables [23, 17].
23 Sparse autoencoders (SAEs) are widely used to extract such features from dense model activations
24 by reconstructing activations using a sparse nonnegative code [3, 12]. This practice implicitly treats
25 SAE latents as candidates for the features used by the model.

26 Recent work shows that this interpretation is fragile. SAEs can learn non-canonical features [20], learn
27 different dictionaries from the same data under different random seeds [26], and exhibit pathologies
28 such as feature absorption [8]. Toy-model recovery studies sharpen the point: recovery is possible
29 only under restrictive sparsity, frequency, or hyperparameter alignment assumptions [9, 6, 11], and
30 high reconstruction quality or decoder recovery need not identify the latent variables [28, 31, 19].

31 The question of this paper is therefore not whether SAEs can sometimes recover features. They can.
32 The question is what the sparse-reconstruction objective incentivizes, even before finite data, encoder
33 amortization, and the optimization landscape enters. The central object of our study is the population
34 one-slot utility

$$\mathcal{U}_\lambda(v; R) = \frac{1}{2} \mathbb{E}[(\langle v, R \rangle - \lambda)_+^2],$$

35 which measures the improvement from adding a nonnegative ℓ_1 -penalized code along direction v
 36 against residual R . Its support-event decomposition shows that the objective scores directions by
 37 thresholded residual mass on latent activation events. Sparse reconstruction is therefore not directly a
 38 feature oracle; it is an event-weighted thresholded residual objective.

39 The analysis separates several objects that are often conflated. A teacher ray may be visible in the
 40 data, but visibility is only raw evidence. The sparse-coding objective then decides whether that
 41 evidence is useful enough above threshold. The SAE encoder then has to realize the corresponding
 42 gate. Finally, even a recovered decoder direction does not always mean that feature’s activation
 43 pattern would also be recovered. Failure modes appear when one object holds but the next one fails:

visibility $\not\Rightarrow$ preference $\not\Rightarrow$ realizability $\not\Rightarrow$ direction recovery $\not\Rightarrow$ support recovery.

44 Each missing implication has a concrete mechanism. Mixed above-threshold events create merging
 45 and hedging-like pressure; event families with varying residual directions create splitting; earlier
 46 surrogate features suppress later evidence, producing absorption-like false negatives; tied or near-tied
 47 event utilities create residual-path alternatives. The following sections establish these points and the
 48 appendix gives the corresponding propositions and proofs.

49 Our contributions are:

- 50 • We identify the one-slot utility and its support-event decomposition as the population objects
 51 controlling one-slot nonnegative sparse reconstruction.
- 52 • We show that one important recovery regime occurs when singleton rank-one events domi-
 53 nate off-singleton contamination, forcing maximizers to localize near teacher directions.
- 54 • We identify mechanisms of the population sparse reconstruction objective that can underlie
 55 several SAE failure modes: merging, splitting, absorption-like false negatives, and residual-
 56 path alternatives can arise before encoder amortization or optimization dynamics.
- 57 • We separate oracle preference from one-ReLU SAE realizability, and decoder direction
 58 recovery from support recovery.
- 59 • We provide a unified interpretation for recent toy-model recovery results, connecting their
 60 observations under a common thresholded utility framework.

61 Sparse reconstruction does not ask “what are the true features?” It asks which directions reduce thresholded
 residual loss. Feature recovery happens only when the data distribution makes true feature directions win
 that scoring game, and when the SAE can realize the winning gate.

62 The paper proceeds as follows. [Section 2](#) defines the positive latent-ray residual model. [Section 3](#)
 63 derives the one-slot utility and support-event decomposition. [Section 4](#) explains recovery and
 64 objective-level failure modes in the oracle setting. [Section 5](#) relates the oracle to one-ReLU SAE
 65 encoders. [Section 6](#) separates direction recovery from support recovery and points to [Section L](#) for
 66 observer-only diagnostics when ground-truth supports are unavailable.

67 2 Exact positive latent-ray model setup

68 We write $[F] := \{1, \dots, F\}$ as a feature index set and $S^{d-1} := \{v \in \mathbb{R}^d : \|v\| = 1\}$ as a unit
 69 sphere. All probabilities and expectations are taken with respect to the population law of the random
 70 activation vector x and the latent coefficients that generate it. We assume finite second moments, for
 71 example $\mathbb{E}\|x\|^2 < \infty$, so that all squared losses and utilities below are finite.

72 We use the linear representation hypothesis, widely adopted in related literature [[17](#), [7](#), [28](#), [6](#), [9](#), [11](#), [18](#)],
 73 as background for our theoretical setup. Throughout, x follows the positive latent-ray model

$$x = \sum_{i=1}^F \alpha_i w_i, \quad \alpha_i \geq 0, \tag{1}$$

74 where $x \in \mathbb{R}^d$ is the observed activation vector, α_i are nonnegative random scalars, $w_i \in \mathbb{R}^d \setminus \{0\}$
 75 are fixed teacher directions, and the rays $\mathbb{R}_+ w_i := \{t w_i : t \geq 0\}$ are distinct. The latent support

76 is $S := \{i : \alpha_i > 0\} \subseteq [F]$. For a context index set $B \subseteq [F]$, chosen or given, and any remaining
 77 subset $T \subseteq [F] \setminus B$, write $E_T(B) := \{S \setminus B = T\}$. Let

$$V_B := \text{span}\{w_i : i \in B\}, \quad P_B : \mathbb{R}^d \rightarrow V_B, \quad \pi_B := I - P_B : \mathbb{R}^d \rightarrow V_B^\perp$$

78 be the orthogonal projections onto the context subspace and its orthogonal complement. Throughout
 79 fixed- B statements are understood on non-vacuous residual charts, i.e. $\dim(V_B) < d$. We define the
 80 projected residual

$$R_B := \pi_B(x) = \sum_{i \notin B} \alpha_i \tilde{w}_i, \quad \tilde{w}_i := \pi_B(w_i).$$

81 On $\{R_B \neq 0\}$, the normalized residual direction is $U_B := R_B / \|R_B\|$.

82 The context B is a way of asking a local question: after some directions have been accounted for,
 83 what residual evidence remains? If the residual on some event points along a single ray, then that
 84 event is clean evidence for a direction. If several support patterns produce the same or nearby residual
 85 direction, the observed activation law can support a direction without uniquely identifying the latent
 86 support semantics behind it.

87 With this residual chart fixed, visibility has a precise meaning. For $j \notin B$ with $\tilde{w}_j \neq 0$, write
 88 $u_j := \tilde{w}_j / \|\tilde{w}_j\|$. We say that the residualized teacher ray $\mathbb{R}_+ \tilde{w}_j$ is *visible at stage B* if the law of
 89 U_B has positive mass at u_j : $\mathbb{P}(U_B = u_j) > 0$. More generally, an observable direction $u \in S^{d-1}$ is
 90 visible at scale ε if $\mathbb{P}(\angle(U_B, u) \leq \varepsilon) > 0$.

91 The clean regime is a rank-one event: the residual points along one ray throughout the event.

92 **Definition 2.1** (Rank-one event). Let $u \in S^{d-1}$ and let $A \subseteq \{R_B \neq 0\}$ be measurable. A is a
rank-one event for the ray $\mathbb{R}_+ u$ if $R_B = \xi u$ with probability one on A for some nonnegative
 random scalar ξ , and $\mathbb{P}(A \cap \{\xi > 0\}) > 0$.

93 Singleton support events $E_{\{j\}}(B)$ are the canonical rank-one events: on $E_{\{j\}}(B)$, the residual is
 94 $R_B = \alpha_j \tilde{w}_j$, so $U_B = u_j$ whenever $\alpha_j > 0$. Thus singleton events make teacher rays directly visible
 95 at the population level unless another support family produces residuals on the same ray; this is
 96 analogous to the separability condition in NMF-style dictionary learning [15].

97 We study a clean positive linear model so that failures cannot be blamed on nonlinearities, signs, finite
 data, or bad optimization. The key observable is the residual direction after removing a context subspace.
 Singleton residual events are the clean evidence for teacher rays.

98 3 What does ℓ_1 sparse coding prefer?

99 For a residual random vector $R \in \mathbb{R}^d$ with $\mathbb{E}\|R\|^2 < \infty$, a unit direction $v \in S^{d-1}$, and sparsity
 100 weight $\lambda \geq 0$, we define the one-slot loss

$$\mathcal{L}(v, a; R) := \mathbb{E} \left[\frac{1}{2} \|R - v a(R)\|_2^2 + \lambda a(R) \right]$$

101 over nonnegative measurable scalar codes $a(R) \geq 0$ with finite second moment. Optimizing this
 102 code gives the oracle utility.

103 This is the first place where the paper's interpretation changes. The one-slot analysis says that a
 104 direction is valuable when enough residual samples project above the sparsity threshold onto it.

105 **Theorem 3.1** (Exact oracle utility). For fixed $v \in S^{d-1}$, the optimal nonnegative one-slot code
 is $a_v^*(R) = (\langle v, R \rangle - \lambda)_+$, and its gain over using no feature is

$$\mathcal{U}_\lambda(v; R) = \frac{1}{2} \mathbb{E}[(\langle v, R \rangle - \lambda)_+^2]. \quad (2)$$

106 The positive part in $(\langle v, R \rangle - \lambda)_+$ says that a sample contributes nothing unless its projection onto
 107 v beats the sparsity threshold. The square then makes large above-threshold residuals especially

108 valuable. Thus sparse coding does not first ask whether v is a teacher atom, it asks whether projection
 109 onto v exceeds the threshold often enough, and with enough magnitude, to pay for the ℓ_1 penalty.
 110 The support-event decomposition then asks where those above-threshold projections come from:
 111 singleton events, mixed events, or broader families of events.

Theorem 3.2 (Support-event decomposition). *The utility decomposes as a sum over residual support events $E_T(B)$. The index $T = \emptyset$ denotes the empty event that contributes no utility:*

$$112 \quad \mathcal{U}_\lambda(v; R_B) = \sum_{\emptyset \neq T \subseteq [F] \setminus B} \frac{1}{2} \mathbb{E} \left[\left(\sum_{i \in T} \alpha_i \langle v, \tilde{w}_i \rangle - \lambda \right)_+^2 \mathbf{1}_{E_T(B)} \right].$$

113 Each term in the sum is the thresholded utility contributed by a single latent support family T . If
 114 $T = \{j\}$, the term is clean singleton evidence for teacher j . If T contains several active atoms, the
 115 term is mixed evidence that may reward a non-teacher direction. This is the main organizing object
 116 of the paper: recovery and non-recovery are comparisons between eventwise utility contributions.

117 The sparse-coding objective has a simple population score. A direction is good when many residuals project above the threshold onto it. The decomposition says which latent support events pay for a direction.

118 4 When does sparse-coding preference recover teacher directions?

119 The support-event decomposition turns feature recovery into an eventwise dominance question. A
 120 teacher direction is preferred in the clean regime where rank-one singleton events dominate all
 121 competing residual events. When this dominance fails, the objective can prefer alternatives to atom
 122 recovery that are analogous to known SAE failure modes in gradient descent training regimes.

123 The theorem below gives the positive case: if clean singleton evidence has a margin over everything
 124 else, the objective localizes near the teacher direction. The following subsections explain what re-
 125 places recovery when that margin is absent. Mixed events create merge pressure, heterogeneous event
 126 families create split pressure, and sequential residual subtraction creates holes and path dependence.

127 4.1 The positive case: rank-one event dominance

128 Let $u_j = \tilde{w}_j / \|\tilde{w}_j\|$. For an angular scale $\varepsilon > 0$, define the singleton dominance gap

$$\mathfrak{M}_{j,\lambda}^B(\varepsilon) := \underbrace{\frac{1}{2} \mathbb{E} \left[\left((\alpha_j \|\tilde{w}_j\| - \lambda)_+^2 - (\alpha_j \|\tilde{w}_j\| \cos \varepsilon - \lambda)_+^2 \right) \mathbf{1}_{E_{\{j\}}(B)} \right]}_{\Delta_{j,\lambda}^B(\varepsilon) \text{ singleton angular margin}} - \underbrace{\frac{1}{2} \mathbb{E} [\|R_B\|^2 \mathbf{1}_{E_{\{j\}}(B)^c}]}_{C_j^B \text{ off-singleton load}}.$$

129 $\Delta_{j,\lambda}^B(\varepsilon)$ is the singleton utility lost by rotating at least ε away from the teacher direction u_j . C_j^B is a
 130 deliberately crude budget for everything outside the singleton event that could compensate for that
 131 loss. If the singleton loss is larger than the off-singleton budget, then no far-away direction can win.

Theorem 4.1 (Singleton dominance implies directional localization). *If $\mathfrak{M}_{j,\lambda}^B(\varepsilon) > 0$, then every global maximizer of $v \mapsto \mathcal{U}_\lambda(v; R_B)$ lies inside the spherical ε -ball around u_j .*

133 Thus ground-truth direction preference is a dominance regime: singleton rank-one evidence must
 134 exceed the total off-singleton competition budget. This gives an eventwise interpretation of extreme-
 135 sparsity recovery results such as [11]: when singleton events dominate, sparse reconstruction can
 136 recover teacher directions. Uniform-margin, low-coherence, support-fragility, random-geometry, and
 137 learned-teacher refinements are collected in Sections F and G.

138 4.2 Merging: one direction compresses several events

139 In the utility field, *merging* means that one direction receives high utility because it compresses several
 140 support events, or several atoms within a mixed event. The sparse-reconstruction objective can then

141 prefer a reconstruction-useful mixed direction over an individual teacher direction. The formal local
142 certificate is a positive tangent derivative away from a teacher direction, stated in [Proposition E.1](#).
143 Singleton-only populations do not strictly favor merges ([Proposition E.2](#)), while independent co-
144 occurrence can favor a merge once the joint projection crosses threshold ([Proposition E.3](#)).

145 The point is not that co-occurrence must be correlated or hierarchical. Correlation is one empirical
146 source of above-threshold mixed events, but the oracle condition is eventwise thresholded transverse
147 residual mass. Put differently, samples that already pay for one direction may also contain residual
148 mass in another direction, so the objective can improve reconstruction by rotating toward a mixture.
149 This is the population version of feature hedging: a learned latent carries a mixture because the
150 mixture improves sparse reconstruction [5].

151 4.3 Splitting: several directions compress one broad family

152 Feature splitting can arise at the objective level when a broad event family does not have a fixed
153 residual direction. No single direction compresses all subevents as well as specialized child directions.
154 [Proposition E.7](#) gives the general event-partition certificate, and [Proposition E.8](#) gives a two-subevent
155 angle-and-threshold specialization.

156 This is distinct from encoder-induced splitting. Objective-level splitting says different subevents
157 prefer different directions even with unrestricted sample-wise gates. Encoder-induced splitting,
158 treated in [Section 5](#), says a single direction may be oracle-preferred while a one-ReLU encoder cannot
159 realize its gate, so several SAE slots may partition the behavior.

160 The nontechnical picture is that a broad feature can be too geometrically heterogeneous for one
161 direction to represent efficiently. The objective then rewards children not because the semantic feature
162 is false, but because its residual evidence arrives through several different rays.

163 4.4 Residual suppression and absorption-like false negatives

164 Merging and splitting are properties of a fixed residual utility field. Absorption-like behavior is
165 different: earlier selections change the residual field. We view absorption as a support-level false-
166 negative pattern where an earlier learned surrogate fires on a subevent, carries part of a broader teacher
167 direction, and suppresses later residual evidence for that direction [8]. The residual-suppression
168 theorem is stated in [Theorem E.4](#); the independent-support construction in [Proposition E.5](#) shows
169 that the same signature can arise without a parent-child generative hierarchy. For an event A , the
170 localized utility used in the formal statement is $\mathcal{U}_{\lambda,A}(v; R) := \frac{1}{2}\mathbb{E}[(\langle v, R \rangle - \lambda)_+^2 \mathbf{1}_A]$. This is the
171 first mechanism where order matters. In a fixed residual field, a teacher direction may have evidence.
172 After a surrogate is selected and subtracted, the same teacher direction may no longer cross threshold
173 on part of its support. The learned feature can then be direction-aligned but have systematic holes in
174 recall.

175 4.5 Residual-path alternatives and seed sensitivity

176 Empirical seed instability refers to SAEs trained on the same data with different random seeds
177 learning different feature sets [26]. We isolate a population residual-path non-identifiability that
178 can serve as an objective-level substrate for this sensitivity. When the oracle utility has multiple
179 tied or near-tied maximizers, residual suppression makes early selections persistent: the selected
180 direction removes its event from the above-threshold residual field, so later utilities depend on the
181 earlier choice. [Proposition E.6](#) constructs an exact population model where tied pair directions give
182 inequivalent non-teacher residual paths.

183 This is not a convergence theorem for SGD. A seed need not create a bad local optimum; in such
184 regimes it can act as a tie-breaker, or near-tie breaker, among residual paths already present in the
185 population objective. Related sparse dictionary learning non-identifiability and multiple-minima
186 phenomena are discussed in [28].

187 [Table 1](#) summarizes the section. Merging and splitting are fixed-field mechanisms; absorption-like
188 holes and residual-path sensitivity require earlier choices because the residual field itself changes.

189 Additional oracle-level material is deferred to the appendix; see [Proposition D.3](#), [Theorem E.9](#),
190 and [section E](#).

Table 1: Mechanism map for the oracle sparse-coding framework. Oracle mechanisms are properties of a fixed residual utility field. Sequential mechanisms arise because earlier learned directions change the residual field.

Mechanism	When it appears	Result
<i>Oracle utility: fixed residual field</i>		
Atom recovery	Singleton rank-one events dominate off-singleton contamination at an angular margin, so oracle maximizers localize near a teacher direction.	Theorem 4.1
Merging	Mixed support events yield larger utility than singleton teacher events; locally, samples above threshold along one direction also contain transverse residual mass.	Proposition E.1 and Proposition E.3
Splitting	A broad event family decomposes into rank-one subevents with different residual directions, and specialized children beat any shared direction.	Proposition E.7 and Proposition E.8
<i>Sequential residual path: earlier choices change later utility</i>		
Residual suppression	A learned surrogate fires on an event supporting another direction and suppresses that event’s above-threshold projection for the original direction.	Theorem E.4 and Proposition E.5
Residual-path sensitivity	Several directions are tied or nearly tied; selecting one removes its associated event from later above-threshold utility fields.	Proposition E.6

191 Together, these cases show the shape of the recovery regime. Teacher directions win when rank-one
 192 events dominate contamination in thresholded utility. Outside that regime, sparse reconstruction can
 193 prefer non-teacher directions that compress mixed events, allocate directions to subevents, or persist
 194 by suppressing residual evidence for the original atom. These mechanisms are population oracle and
 195 residual-greedy mechanisms, before encoder amortization or optimization dynamics.

196 Teacher features win only when their clean rank-one events dominate the other support events. If mixed or heterogeneous events carry more above-threshold residual mass, the same objective prefers merged, split, or path-dependent alternatives.

197 5 Does an SAE realize the oracle preference?

198 The oracle in [Section 3](#) was allowed to choose a fresh nonnegative code for every sample. A SAE
 199 slot is more constrained: its activation must be produced by one encoder row. This adds a new layer
 200 to the implication chain. A direction can be preferred by sparse coding, yet fail to be realized by a
 201 one-ReLU SAE gate.

202 Let K be the number of SAE slots. Let $D = [d_1, \dots, d_K] \in \mathbb{R}^{d \times K}$ have unit-norm columns, let
 203 $E \in \mathbb{R}^{K \times d}$ be the encoder matrix, let $b_{\text{enc}} \in \mathbb{R}^K$ and $b_{\text{dec}} \in \mathbb{R}^d$ be encoder and decoder biases, and
 204 let $\sigma(t) = t_+$ act coordinatewise. The encoder is *untied*: E is not constrained to be equal to D^\top .
 205 The encoder input is $\nu = x - b_{\text{dec}}$, the code is $z = \sigma(E\nu + b_{\text{enc}}) \in \mathbb{R}_+^K$, and the reconstruction is
 206 $\hat{x} = Dz + b_{\text{dec}}$. For slot r , write $R_{-r} := \nu - \sum_{s \neq r} d_s z_s$ and $y_r := d_r^\top R_{-r} - \lambda$.

207 The blockwise algebra is simple but important. Once the other slots and the decoder direction are
 208 fixed, the encoder row is just trying to imitate the oracle threshold target $y_r = d_r^\top R_{-r} - \lambda$, with a
 209 nonnegative one-ReLU output.

Theorem 5.1 (Block decomposition). *For a fixed unit decoder direction, fixed other slots, and a fixed sample value of the leave-one-out residual, the slot loss satisfies the pointwise identity*

$$210 \quad \frac{1}{2} \|R_{-r} - d_r z_r\|^2 + \lambda z_r = \frac{1}{2} \|R_{-r}\|^2 - \frac{1}{2} y_r^2 + \frac{1}{2} (z_r - y_r)^2.$$

Under this conditioning, optimizing one encoder row is squared-error regression onto y_r with a nonnegative one-ReLU output; the unrestricted pointwise nonnegative optimum is $(y_r)_+$.

211 Thus the sample-wise sparse-coding gate for direction d and residual R is $g_d^* = (d^\top R - \lambda)_+$.
 212 The SAE row can match this oracle only if the same random function can be written as a single
 213 affine-ReLU of the encoder input. To state this precisely, define

$$\mathcal{A}_\nu := \{ (e^\top \nu + b)_+ : e \in \mathbb{R}^d, b \in \mathbb{R} \} \subset L_+^2(\Omega)$$

214 and, for square-integrable random variable gate $z \in L_+^2$,

$$\Gamma_\lambda(d, z; R) := \mathbb{E} \left[z(d^\top R - \lambda) - \frac{1}{2} z^2 \right], \quad \mathcal{A}_\lambda(d; R, \nu) := \sup_{z \in \mathcal{A}_\nu} \Gamma_\lambda(d, z; R).$$

215 Here \mathcal{A}_λ is the one-row SAE version of the oracle utility: same decoder direction, same residual, but
 216 the gate must be a single affine-ReLU of ν . The full envelope identity is deferred to [Theorem H.2](#);
 217 the main consequence is the closure criterion below.

Theorem 5.2 (Exact one-ReLU realizability and closure criterion). *Assume $R, \nu \in L^2$, and fix $d \in S^{d-1}$. Let $\overline{\mathcal{A}_\nu}^{L^2}$ denote the closure of \mathcal{A}_ν in $L^2(\Omega)$. Then*

218
$$\mathcal{A}_\lambda(d; R, \nu) = \mathcal{U}_\lambda(d; R) \iff g_d^* \in \overline{\mathcal{A}_\nu}^{L^2}.$$

If the supremum defining $\mathcal{A}_\lambda(d; R, \nu)$ is attained by some $z_{e,b} = (e^\top \nu + b)_+$, then equality holds if and only if $(e^\top \nu + b)_+ = g_d^$ a.s. Thus exact one-ReLU realizability is the attained case of the L^2 -closure criterion.*

219 This is the encoder-realizability gap. If $g_d^* \notin \overline{\mathcal{A}_\nu}^{L^2}$, the one-row SAE value is strictly below the
 220 sample-wise oracle value even when its decoder direction is correct. If g_d^* lies in the closure but is not
 221 exactly affine-ReLU realizable, the oracle value can be approached but is not attained by a single row.
 222 Several slots may then partition a single oracle gate and reduce this gap, producing an architectural
 223 splitting mechanism distinct from the event-level splitting in [Section 4.3](#).

224 [Appendix M](#) gives a small sanity check for the realizability layer. We first train SAE1 on the ground-
 225 truth generator, then train SAE2 either directly on SAE1's code or on SAE1's reconstructed hidden
 226 states. Direct code training is the easy endpoint: SAE2 receives a nonnegative coordinate repre-
 227 sentation and recovers SAE1 features with support F1 0.937, compared with 0.583 for GT→SAE1
 228 and 0.600 for GT→SAE2. Training on reconstructed hidden states is harder, but the representation-
 229 relative pattern remains: SAE1→SAE2 support F1 is 0.804 with learned decoders and 0.592 in the
 230 fixed-decoder output-student ablation, compared to the GT→SAE1 values 0.583 and 0.571.

231 5.1 Transfer from latent residuals to SAE slots

232 There is one more transfer step. The latent sparse-coding score uses $R_B = \pi_B x$, while slot r sees
 233 R_{-r} . Define the perturbation scale $\eta(R, S) := \frac{1}{2} \|R - S\|_{L^2} (\|R\|_{L^2} + \|S\|_{L^2})$, with the full bound
 234 stated in [Lemma I.2](#). Let $\eta_{B,-r} := \eta(R_B, R_{-r})$ and define the amortization gap

$$\text{Gap}_\lambda(u; R_{-r}, \nu) := \mathcal{U}_\lambda(u; R_{-r}) - \mathcal{A}_\lambda(u; R_{-r}, \nu),$$

235 the utility lost at direction u because the gate must be a single affine-ReLU of ν . For a competitor set
 236 \mathcal{C} , the corresponding amortized SAE margin is

$$\text{Marg}_\mathcal{C}^{\text{SAE}}(u) := \mathcal{A}_\lambda(u; R_{-r}, \nu) - \sup_{v \in \mathcal{C}} \mathcal{A}_\lambda(v; R_{-r}, \nu).$$

237 The latent margin γ is the amount by which a direction beats its competitors in the clean residual field
 238 R_B . Two costs spend that margin before it reaches an SAE: residual mismatch and gate amortization.

Theorem 5.3 (Residual-field transfer). *Let $\mathcal{C} \subseteq S^{d-1}$ be a set of competitor directions, and suppose $u \in S^{d-1}$ has latent margin $\gamma := \mathcal{U}_\lambda(u; R_B) - \sup_{v \in \mathcal{C}} \mathcal{U}_\lambda(v; R_B) > 0$. Then $\text{Marg}_\mathcal{C}^{\text{SAE}}(u) \geq \gamma - 2\eta_{B,-r} - \text{Gap}_\lambda(u; R_{-r}, \nu)$. Hence the ranking transfers whenever this lower bound is positive.*

240 The factor $2\eta_{B,-r}$ is the cost of changing residual fields; the gap term is the cost of restricting the
 241 gate to one affine-ReLU. Exact residual compatibility is

$$R_{-r} = R_B \iff b_{\text{dec}} + \sum_{s \neq r} d_s z_s = P_B x \quad \text{almost surely.}$$

242 A decoder bias inside V_B leaves the projected residual component unchanged, but a nonzero projected
 243 bias shifts thresholds and can create utility on empty residual-support events. Amortization can
 244 reverse an oracle ranking only through direction-dependent gaps; affine residual regimes eliminate
 245 those gaps, while pair-only worlds provide stress tests where both oracle and exact sparse coding
 246 prefer non-ground-truth directions (Sections H.7 and I).

247 Even if the oracle prefers a direction, an SAE slot still has to realize the right gate using a one-ReLU
 encoder row and the leave-one-out residual. Oracle preference and SAE realization are separate layers.

248 6 Do recovered directions imply support semantics?

249 The previous section separated oracle preference from encoder realizability. The next separation
 250 is semantic: even if a decoder direction is recovered, this need not mean that its gate fires on the
 251 intended latent examples. Direction recovery and support recovery are different objects.

252 The reason is visible in one moment equation. Let $a(x) \geq 0$ be an activation such that $\mathbb{E}[\|R_B\|a] <$
 253 ∞ and $\mathbb{E}[\alpha_i a] < \infty$ for every $i \notin B$. Define $\mu(a; B) := \mathbb{E}[R_B a] = \sum_{i \notin B} \mathbb{E}[\alpha_i a] \tilde{w}_i$.

254 When the decoder is fit by least squares to the gate, and in the corresponding stationary condition
 255 for the SAE decoder, the decoder direction aligns with this gated residual moment. The decoder
 256 therefore sees an average vector weighted by the gate. If the gate also fires when other atoms are
 257 active, their residual directions enter the average; if it fires on only a tiny but geometrically pure
 258 subset, the direction can look correct while support recall is poor.

259 **Theorem 6.1** (Moment characterization of direction alignment). *For $j \notin B$ with $\tilde{w}_j \neq 0$, let
 $a \geq 0$ satisfy $\mathbb{E}[\|R_B\|a] < \infty$ and $\mathbb{E}[\alpha_i a] < \infty$ for every $i \notin B$, and suppose $\mu(a; B) \neq 0$.
 Then $\mu(a; B) \in \text{span}(\tilde{w}_j)$ iff the off- j gated moment lies in $\text{span}(\tilde{w}_j)$. If the residualized
 atoms are linearly independent, this is equivalent to $\alpha_i a = 0$ a.s. for every $i \notin B, i \neq j$.*

260 Under linear independence, exact nonzero decoder alignment with a teacher ray is equivalent to zero
 261 gated co-activation with every other residualized atom. This is stronger than “the gate fires when
 262 the atom is present.” It is also different from support recovery. The minimal counterexamples in
 263 Propositions J.1 and J.2 show both directions of failure: a gate can recover direction with arbitrarily
 264 poor support recall, and a support-faithful gate can yield a rotated decoder direction.

265 The claim that direction recovery implies support recovery is also disproven in [31]; here we isolate
 266 the population event-geometry behind this separation. Below, we show that the two objects coincide
 267 in the clean regime already used for oracle recovery: rank-one event structure.

268 **Theorem 6.2** (Coincidence on rank-one events). *If A is a rank-one event with $R_B = \xi u$ and
 the gate $a \geq 0$ is supported on A with $0 < \mathbb{E}[\xi a \mathbf{1}_A] < \infty$, then $\mu(a; B)$ is parallel to u .*

269 This is only a sufficient coincidence regime: event support plus positive gated amplitude on a rank-one
 270 residual event yields direction alignment. It does not say that arbitrary gates aligned with u recover A .
 271 The robust moment bound and support-fragility refinements are in Section J; learned-teacher simplex
 272 symmetry gives another special coincidence regime by cancelling exchangeable off-target moments
 273 (Section G). These are structured exceptions, not generic consequences of sparse reconstruction.

274 6.1 Sparsity matching is not support-event matching

275 Recent toy-model work emphasizes sparsity-level and firing-frequency alignment as recovery-relevant
 276 quantities [6, 9]. In our framework these quantities are real necessary shadows of support recovery,

277 but they are not the support-event law itself. Let $Z_i := \mathbf{1}_{\{\alpha_i > 0\}}$, $S := \{i : Z_i = 1\}$, $N := |S|$,
 278 $p_i := \mathbb{P}(Z_i = 1)$, and $s_* := \mathbb{E}N$. Exact support recovery implies samplewise sparsity recovery,
 279 which implies expected sparsity matching: $\widehat{S} = S$ a.s. $\implies |\widehat{S}| = |S|$ a.s. $\implies \mathbb{E}|\widehat{S}| = s_*$. The
 280 converses fail because expected sparsity fixes only $\sum_i p_i$, and marginal-frequency matching fixes
 281 only the one-feature marginals. The sparse-reconstruction objective decomposes over full support
 282 events $E_T(B)$, so it depends on joint probabilities $\mathbb{P}(S = T)$, eventwise amplitudes, and residual
 283 geometry. The missing information already appears at second order:

$$\text{Var}(N) = \sum_i p_i(1 - p_i) + 2 \sum_{i < j} \text{Cov}(Z_i, Z_j), \quad \mathbb{E} \binom{N}{2} = \sum_{i < j} \mathbb{P}(Z_i = Z_j = 1).$$

284 Pairwise co-activation events are exactly the lowest-order mixed support events in [Theorem 3.2](#).
 285 They can be invisible to a global sparsity target while still dominating the thresholded utility. Thus
 286 sparsity and frequency matching can remove one mismatch between teacher and learned codes,
 287 but they do not replace support-event dominance. [Section K](#) gives the exact two-feature marginal-
 288 frequency separation, amplitude-threshold mismatch under the correct orthogonal dictionary, power-
 289 law asymptotics, load transitions, and pairwise-only support constructions.

290 [Appendix M](#) also illustrates the direction/support distinction. In fixed-decoder ablations, decoder
 291 rows are correct by construction. Nevertheless, GT support recovery remains poor: direction-matched
 292 GT→SAE2 support F1 is 0.569 when SAE2 is trained on SAE1’s code, and 0.391 when SAE2 is
 293 trained on SAE1 reconstructed hidden states with the GT decoder fixed. In the latter case, L_0 is also
 294 closely matched to the teacher target (1.153 versus 1.175) and explained variance is high (0.968),
 295 so the failure is not just a gross sparsity or reconstruction mismatch. The closest support-matched
 296 features are geometrically nearby, about 0.94; the issue is that their support events do not match.
 297 Also, while GT supports are independent, learned supports have substantial off-diagonal correlations.

298 The observer-level version is even weaker. In real activations the latent support S is not observed,
 299 so hidden states alone cannot identify latent support semantics. What can be certified is residual
 300 geometry: observable residual-direction atoms or caps, directional preference of the observed sparse-
 301 coding utility field, and whether learned features match those caps by direction, precision, recall,
 302 intrinsic ambiguity, and gate consistency. [Section L](#) gives the formal observer-only diagnostics.

303 A decoder direction and a feature’s firing pattern are different semantic objects. Direction alignment
 is a gated-moment condition; support recovery is an event condition. Sparsity and frequency matching
 preserve useful low-order statistics, but the objective depends on the full support-event law.

304 7 Conclusion

305 This paper isolates what nonnegative sparse reconstruction prefers before an SAE encoder, finite
 306 data, or gradient dynamics enter. The one-slot utility $\mathcal{U}_\lambda(v; R) = \frac{1}{2} \mathbb{E}[(\langle v, R \rangle - \lambda)_+^2]$ is a population
 307 thresholded gain, and its support-event decomposition shows which latent events make a direction
 308 valuable. Teacher directions win only in structured dominance regimes, where rank-one or single-
 309 feature events dominate off-target contamination above threshold. When this dominance fails, the
 310 same objective can prefer non-canonical but reconstruction-useful directions: mixed events create
 311 merging, heterogeneous event families create splitting, earlier surrogates create absorption-like false
 312 negatives, and near-tied event utilities create residual-path alternatives.

313 The claims fit together through the implication chain

$$\text{visibility} \not\Rightarrow \text{preference} \not\Rightarrow \text{realizability} \not\Rightarrow \text{direction recovery} \not\Rightarrow \text{support recovery}.$$

314 Each missing implication is witnessed by an explicit construction in [Section C](#). Thus sparse recon-
 315 struction can recover ground-truth features, but only when the data distribution, threshold, geometry,
 316 and encoder class align; outside those regimes, such failures should not be attributed only to bad
 317 optimization or architecture, because population objectives can already incentivize them.

318 The limitations are structural. The nonnegative cone in [Equation \(1\)](#) is essential for the residual-ray
 319 and support-semantics interpretation of [Theorem 3.2](#). Signed activations, nonlinear mixing, SAE
 320 optimization dynamics and prediction of SAE preferences, finite-sample observer diagnostics, and
 321 empirical validation in pretrained-model SAEs remain open. The fixed- B analysis applies only on
 322 non-vacuous residual charts with $\dim(V_B) < d$, and target-feature claims require $\pi_B w_j \neq 0$.

References

- 323
- 324 [1] M. Aharon, M. Elad, and A. Bruckstein. K-svd: An algorithm for designing overcomplete
325 dictionaries for sparse representation. *IEEE Transactions on Signal Processing*, 54(11):4311–
326 4322, 2006.
- 327 [2] Sanjeev Arora, Rong Ge, Tengyu Ma, and Ankur Moitra. Simple, efficient, and neural algorithms
328 for sparse coding. In Peter Grünwald, Elad Hazan, and Satyen Kale, editors, *Proceedings of
329 The 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine Learning
330 Research*, pages 113–149, Paris, France, 03–06 Jul 2015. PMLR.
- 331 [3] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Con-
332 erly, Nick Turner, Cem Anil, Carson Denison, Amanda Aspell, Robert Lasenby, Yifan Wu,
333 Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex
334 Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter,
335 Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language
336 models with dictionary learning. *Transformer Circuits Thread*, 2023. [https://transformer-
337 circuits.pub/2023/monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 338 [4] E.J. Candes and T. Tao. Decoding by linear programming. *IEEE Transactions on Information
339 Theory*, 51(12):4203–4215, 2005.
- 340 [5] David Chanin, Tomáš Dulka, and Adrià Garriga-Alonso. Feature hedging: Correlated features
341 break narrow sparse autoencoders, 2025.
- 342 [6] David Chanin and Adrià Garriga-Alonso. Sparse but wrong: Incorrect l0 leads to incorrect
343 features in sparse autoencoders, 2025.
- 344 [7] David Chanin and Adrià Garriga-Alonso. Synthsaebench: Evaluating sparse autoencoders on
345 scalable realistic synthetic data, 2026.
- 346 [8] David Chanin, James Wilken-Smith, Tomáš Dulka, Hardik Bhatnagar, Satvik Golechha, and
347 Joseph Isaac Bloom. A is for absorption: Studying feature splitting and absorption in sparse
348 autoencoders. In *The Thirty-ninth Annual Conference on Neural Information Processing
349 Systems*, 2026.
- 350 [9] Siyu Chen, Heejune Sheen, Xuyuan Xiong, Tianhao Wang, and Zhuoran Yang. Taming
351 polysemanticity in LLMs: Theory-grounded feature recovery via sparse autoencoders. In *The
352 Fourteenth International Conference on Learning Representations*, 2026.
- 353 [10] Valérie Costa, Thomas Fel, Ekdeep Singh Lubana, Bahareh Tolooshams, and Demba Ba.
354 Evaluating sparse autoencoders: From shallow design to matching pursuit, 2025.
- 355 [11] Jingyi Cui, Qi Zhang, Yifei Wang, and Yisen Wang. On the limits of sparse autoencoders: A
356 theoretical framework and reweighted remedy. In *The Fourteenth International Conference on
357 Learning Representations*, 2026.
- 358 [12] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse
359 autoencoders find highly interpretable features in language models, 2023.
- 360 [13] Jean-François Determe, Jérôme Louveaux, Laurent Jacques, and François Horlin. On the exact
361 recovery condition of simultaneous orthogonal matching pursuit. *IEEE Signal Processing
362 Letters*, 23(1):164–168, 2016.
- 363 [14] David L. Donoho and Michael Elad. Optimally sparse representation in general (nonorthogonal)
364 dictionaries via ℓ^1 minimization. *Proceedings of the National Academy of Sciences*, 100(5):2197–
365 2202, 2003.
- 366 [15] David L. Donoho and Victoria Stodden. When does non-negative matrix factorization give a
367 correct decomposition into parts? In *NIPS*, pages 1141–1148, 2003.
- 368 [16] D.L. Donoho, M. Elad, and V.N. Temlyakov. Stable recovery of sparse overcomplete repre-
369 sentations in the presence of noise. *IEEE Transactions on Information Theory*, 52(1):6–18,
370 2006.

- 371 [17] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna
372 Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam
373 McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy
374 models of superposition. *Transformer Circuits Thread*, 2022.
- 375 [18] Nikhil Garg, Jon Kleinberg, and Kenny Peng. How many features can a language model store
376 under the linear representation hypothesis?, 2026.
- 377 [19] Anton Korznikov, Andrey Galichin, Alexey Dontsov, Oleg Rogov, Ivan Oseledets, and Elena
378 Tutubalina. Sanity checks for sparse autoencoders: Do saes beat random baselines?, 2026.
- 379 [20] Patrick Leask, Bart Bussmann, Michael Pearce, Joseph Bloom, Curt Tigges, Noura Al
380 Moubayed, Lee Sharkey, and Neel Nanda. Sparse autoencoders do not find canonical units of
381 analysis, 2025.
- 382 [21] Julien Mairal, Francis R. Bach, Jean Ponce, and Guillermo Sapiro. Online dictionary learning
383 for sparse coding. In *ICML*, pages 689–696, 2009.
- 384 [22] S.G. Mallat and Zhifeng Zhang. Matching pursuits with time-frequency dictionaries. *IEEE*
385 *Transactions on Signal Processing*, 41(12):3397–3415, 1993.
- 386 [23] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.
387 Zoom in: An introduction to circuits. *Distill*, 2020. <https://distill.pub/2020/circuits/zoom-in>.
- 388 [24] Bruno A. Olshausen and David J. Field. Sparse coding with an overcomplete basis set: A
389 strategy employed by v1? *Vision Research*, 37(23):3311–3325, 1997.
- 390 [25] Charles O’Neill, Alim Gumran, and David Klindt. Compute optimal inference and provable
391 amortisation gap in sparse autoencoders. In *Forty-second International Conference on Machine*
392 *Learning*, 2025.
- 393 [26] Gonçalo Paulo and Nora Belrose. Sparse autoencoders trained on the same data learn different
394 features, 2025.
- 395 [27] Daniel A. Spielman, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries.
396 In Shie Mannor, Nathan Srebro, and Robert C. Williamson, editors, *Proceedings of the 25th*
397 *Annual Conference on Learning Theory*, volume 23 of *Proceedings of Machine Learning*
398 *Research*, pages 37.1–37.18, Edinburgh, Scotland, 25–27 Jun 2012. PMLR.
- 399 [28] Yiming Tang, Harshvardhan Saini, Zhaoqian Yao, Zheng Lin, Yizhen Liao, Jingyi Cui, Yisen
400 Wang, Mengnan Du, and Dianbo Liu. A unified theory of sparse dictionary learning in
401 mechanistic interpretability: Piecewise biconvexity and spurious minima, 2026.
- 402 [29] Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen,
403 Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L
404 Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers,
405 Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan.
406 Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer*
407 *Circuits Thread*, 2024.
- 408 [30] J.A. Tropp. Greed is good: algorithmic results for sparse approximation. *IEEE Transactions on*
409 *Information Theory*, 50(10):2231–2242, 2004.
- 410 [31] Zhenyu Zhu, Marco Fumero, Francesco Locatello, and Volkan Cevher. Diagnosing and fixing
411 latent recovery in sparse autoencoders. In *ICLR 2026 Workshop on Unifying Concept*
412 *Representation Learning*, 2026.

413 A Appendix roadmap

414 The appendix collects the formal material behind the main claims. It preserves the deferred theorem
415 and proposition statements, gives complete derivations, includes refinements that clarify positive
416 regimes, and gives stress tests for the boundary between oracle preference and amortized realization.

417 **Section C** gives the counterexamples behind the implication chain. **Sections D to F** contain the
418 oracle derivations, mechanisms, and positive recovery refinements. **Section G** studies structured
419 learned-teacher geometries. **Sections H and I** give the SAE bridge and residual-transfer details.
420 **Sections J and K** contain the direction-support and sparsity-frequency branches. **Section L** contains
421 observer-only audit refinements. **Section M** gives synthetic two-stage SAE illustrations for the
422 realizability layer and the direction/support separation.

423 B Related work

424 The related literature gives the recovery side of the story and the pathology side of the story. The
425 point of comparison here is the population sparse-reconstruction objective: when the assumptions
426 behind recovery guarantees weaken, which residual events does that objective reward?

427

Prior recovery work asks when sparse methods recover features. The complementary question here is
what the population sparse-reconstruction objective rewards when recovery conditions do not dominate.

428 **Sparse recovery and dictionary learning.** Classical sparse recovery studies when sparse codes or
429 dictionaries are identifiable from data, typically under incoherence, restricted-isometry, separability,
430 or sparse-use assumptions [14, 4, 27, 2]. Classical dictionary-learning algorithms such as sparse
431 coding, K-SVD, and online dictionary learning optimize reconstruction with sparsity constraints and
432 alternate between code inference and dictionary updates [24, 1, 21]. This literature primarily gives
433 sufficient conditions for recovery. Our question is complementary: when recovery conditions do not
434 dominate the population objective, which residual support events does sparse reconstruction reward,
435 and what non-teacher directions can become optimal?

436 **Matching pursuit and greedy sparse approximation.** The one-slot oracle in **Section 3** is the pop-
437 ulation analogue of a nonnegative, ℓ_1 -thresholded matching-pursuit step. Classical matching pursuit
438 selects atoms by greedy residual reconstruction gain [22]; OMP and related analyses characterize
439 when greedy residual correlations recover the correct support under dictionary assumptions [30, 16].
440 Simultaneous matching pursuit extends this logic to multiple signals and common-support recovery
441 [13]. Our setting differs in two ways. First, the direction ranges over the sphere rather than over a
442 fixed dictionary. Second, latent supports vary across samples, so the relevant object is not a single
443 support to recover but the support-event decomposition of the population gain. This decomposi-
444 tion identifies when greedy sparse reconstruction favors teacher directions, merged directions, split
445 subevents, absorption-like false negatives, or residual-path alternatives.

446 **Sparse autoencoders and feature recovery.** SAEs are widely used to extract interpretable features
447 from neural activations [12, 3, 29]. Recent theory studies when such features recover known ground
448 truth in toy latent models. Recovery can hold under restrictive sparsity, frequency, or hyperparameter
449 alignment assumptions [9, 6, 11, 7]. These results align with our positive regimes: extreme sparsity,
450 matched firing frequencies, and correct sparsity levels are ways of making thresholded support events
451 identify teacher directions. Our contribution is not another recovery guarantee for a specific SAE
452 algorithm; it is an event-level explanation of why the same sparse-reconstruction objective can either
453 support recovery or prefer non-canonical directions.

454 **Non-canonical features and SAE pathologies.** Empirical and theoretical work has documented
455 non-canonical SAE features, feature absorption, feature hedging, seed-dependent feature sets, and
456 other failures [20, 8, 5, 26]. We do not claim priority for observing these phenomena. We isolate
457 objective-level mechanisms for a central subset of them. Mixed above-threshold events create
458 merge and hedging-like pressure; variable rank-one subevents create event-partition splitting; earlier
459 surrogates create absorption-like false negatives by suppressing residual evidence; and tied or near-

460 tied event utilities create residual-path non-identifiability. These mechanisms occur before finite
461 samples, encoder amortization, gradient dynamics, or local minima.

462 **Direction recovery versus code or support recovery.** Several recent works emphasize that recon-
463 struction quality or dictionary recovery does not imply recovery of latent variables [28, 31]. Zhu et
464 al. [31] are especially close to our direction/support separation: they show that decoder recovery
465 does not guarantee latent-code recovery. The work [25] also argues that standard SAE encoders
466 cannot optimally represent the underlying code, and that expressivity of encoder class is important
467 for recovery. We refine this distinction for interpretability semantics. Decoder direction recovery is
468 a gated-moment condition, while support recovery is an event-level condition; neither implies the
469 other in general. They coincide only under additional structure, such as rank-one event support or
470 symmetric cancellation in learned geometries.

471 **Matching-pursuit SAEs and broader SDL theory.** Recent MP-SAE work uses residual-guided
472 matching-pursuit inference inside SAE architectures to better handle correlated, hierarchical, or
473 conditionally orthogonal features [10]. This is architecturally adjacent to our oracle, but our use
474 is diagnostic rather than constructive: residual-guided sparse reconstruction can itself prefer non-
475 canonical directions at the population level. This also complements unified SDL theory, which studies
476 optimization landscapes, non-identifiability, and spurious minima for sparse dictionary learning
477 methods [28]. Our analysis is orthogonal: we study the population thresholded sparse-coding gain
478 first, then separately track encoder representability and direction/support semantics.

479 C Separation witnesses

480 The implication chain in the introduction needs exact witnesses. Each construction below shows that
481 one semantic object can hold without the next one in the chain, so none of the missing implications is
482 silently assumed.

483 Each witness breaks one tempting implication: visibility, preference, realizability, direction recovery, and
support recovery are different objects.

484 **Proposition C.1** (Separation witnesses). *The implication chain separated in the introduction
cannot be collapsed: visibility need not imply oracle preference; oracle preference need not
imply one-ReLU realizability; one-ReLU realizability need not imply teacher-direction recovery;
teacher-direction recovery need not imply support recovery; and support recovery need not
imply teacher-direction recovery.*

485 *Proof.* Visibility need not imply preference: in Proposition E.3, with $p \in (0, 1)$ and $A \leq \lambda < \sqrt{2}A$,
486 the singleton event for u has positive probability, so u is visible, but $\mathcal{U}_\lambda(u; x) = 0 < \mathcal{U}_\lambda((u +$
487 $c)/\sqrt{2}; x)$.

488 Preference need not imply one-ReLU realizability with respect to the encoder input: take the encoder
489 input to be a scalar random variable ν with $\mathbb{E}[\nu^4] < \infty$ and support containing an interval, set
490 $R(\nu) = \nu^2 u$, and take $\lambda = 0$.

491 One-ReLU realizability need not imply teacher-direction recovery: in the same independent two-atom
492 affine model used in Proposition E.3, the merged direction $m = (u + c)/\sqrt{2}$ has affine-ReLU gate
493 $(m^\top x - \lambda)_+$ and is oracle-preferred in the threshold regime, but m is not a teacher ray.

494 Teacher-direction recovery without support recovery is Proposition J.1. Support recovery without
495 teacher-direction recovery is Proposition J.2. \square

496 D Core oracle derivations

497 These derivations support Sections 2–3. They establish what the oracle utility measures, why residual
498 support events decompose its utility, and how marginal insertion relates the one-slot calculation to
499 sparse-coding improvement.

500

The central algebra is: singleton events create residual-ray atoms, one-slot sparse coding gives the thresholded utility, and support events split that utility into interpretable pieces.

Proposition D.1 (Residual-ray observability from singleton witnesses). *Fix $j \notin B$. Assume*

$$\tilde{w}_j \neq 0, \quad \mathbb{P}(S \setminus B = \{j\}) > 0.$$

Then the law of $U_B := R_B / \|R_B\|$ on $\{R_B \neq 0\}$ has an atom at

$$\tilde{u}_j := \frac{\tilde{w}_j}{\|\tilde{w}_j\|}.$$

501

Moreover, suppose that for every family $T \subseteq [F] \setminus B$ with $\mathbb{P}(E_T(B)) > 0$, the following hold:

$$j \notin T, T \neq \emptyset \implies \mathbb{R}_+ \tilde{w}_j \cap \text{cone}_+ \{\tilde{w}_i : i \in T\} = \{0\},$$

and

$$j \in T, T \neq \{j\} \implies \mathbb{R} \tilde{w}_j \cap \text{cone}_+ \{\tilde{w}_i : i \in T \setminus \{j\}\} = \{0\}.$$

Then no other positive residual support family generates the same ray, so the atom is semantically unambiguous relative to the assumed latent factorization.

502 D.1 Proof of Proposition D.1

503 On the singleton event $E_{\{j\}}(B)$, one has

$$R_B = \alpha_j \tilde{w}_j,$$

504 so on $E_{\{j\}}(B) \cap \{\alpha_j > 0\}$,

$$U_B = \frac{\tilde{w}_j}{\|\tilde{w}_j\|} = \tilde{u}_j.$$

505 Thus the law of U_B assigns at least the positive mass $\mathbb{P}(E_{\{j\}}(B) \cap \{\alpha_j > 0\})$ to the singleton
506 direction \tilde{u}_j . This is residual-ray observability, not latent-support identifiability from x alone. The
507 cone-separation condition

$$\mathbb{R}_+ \tilde{w}_j \cap \text{cone}_+ \{\tilde{w}_i : i \in T\} = \{0\}$$

508 for every other support family T rules out semantic collisions with other positive support families
509 under the assumed latent factorization.

510 For $T \not\ni j$, the first cone condition rules out $R_B \in \mathbb{R}_+ \tilde{w}_j$ on $E_T(B)$. For $T \ni j, T \neq \{j\}$, write

$$R_B = \alpha_j \tilde{w}_j + \sum_{i \in T \setminus \{j\}} \alpha_i \tilde{w}_i.$$

511 If $R_B \in \mathbb{R}_+ \tilde{w}_j$, then

$$\sum_{i \in T \setminus \{j\}} \alpha_i \tilde{w}_i \in \mathbb{R} \tilde{w}_j,$$

512 contradicting the second cone condition.

513 D.2 Proof of Theorem 3.1

514 For fixed $R = r$,

$$\frac{1}{2} \|r - va\|_2^2 + \lambda a = \frac{1}{2} \|r\|_2^2 - a \langle v, r \rangle + \frac{1}{2} a^2 + \lambda a.$$

515 This is a strictly convex quadratic in $a \geq 0$. The unconstrained minimizer is $\langle v, r \rangle - \lambda$, so the
516 constrained minimizer is $(\langle v, r \rangle - \lambda)_+$. Substituting back yields the gain

$$\frac{1}{2} (\langle v, r \rangle - \lambda)_+^2.$$

517 Taking expectations gives the stated utility.

518 **D.3 Proof of Theorem 3.2**

519 The support events $\{E_T(B)\}_{T \subseteq [F] \setminus B}$ form a measurable partition. On $E_T(B)$,

$$R_B = \sum_{i \in T} \alpha_i \tilde{w}_i.$$

520 Substitute this into Theorem 3.1 and sum over the partition. Thus

$$\mathcal{U}_{\lambda, T}^B(v) := \frac{1}{2} \mathbb{E} \left[\left(\sum_{i \in T} \alpha_i \langle v, \tilde{w}_i \rangle - \lambda \right)_+^2 \mathbf{1}_{E_T(B)} \right],$$

521 and the empty event contributes zero for the projected residual R_B .

Proposition D.2 (utility gradient). *Let $R \in L^2(\mathbb{R}^d)$. Then*

$$\nabla_v \mathcal{U}_\lambda(v; R) = \mathbb{E}[(\langle v, R \rangle - \lambda)_+ R],$$

522 *and the Riemannian gradient on the sphere is*

$$\text{grad}_{S^{d-1}} \mathcal{U}_\lambda(v; R) = P_{v^\perp} \mathbb{E}[(\langle v, R \rangle - \lambda)_+ R].$$

523 **D.4 Proof of Proposition D.2**

524 Let $\phi(t) = \frac{1}{2}(t - \lambda)_+^2$, so $\phi'(t) = (t - \lambda)_+$. Differentiate under the expectation:

$$D\mathcal{U}_\lambda(v)[\delta] = \mathbb{E}[\phi'(\langle v, R \rangle) \langle \delta, R \rangle] = \langle \mathbb{E}[(\langle v, R \rangle - \lambda)_+ R], \delta \rangle.$$

525 This gives the Euclidean gradient. Project onto v^\perp to obtain the sphere gradient.

Proposition D.3 (marginal insertion bound). *Let*

$$\Phi_\lambda(D) := \mathbb{E} \left[\min_{z \geq 0} \frac{1}{2} \|x - Dz\|_2^2 + \lambda \mathbf{1}^\top z \right].$$

526 *Fix a $(K - 1)$ -slot dictionary D_{-r} and let R_{-r}^* be the corresponding leave-one-out exact residual. Then for every unit vector v ,*

$$\Phi_\lambda(D_{-r}) - \Phi_\lambda(D_{-r} \cup \{v\}) \geq \mathcal{U}_\lambda(v; R_{-r}^*).$$

527 **D.5 Proof of Proposition D.3**

528 For each sample x , start from the exact $(K - 1)$ -slot optimum and optimize only the new coefficient
529 $a \geq 0$ along v . By Theorem 3.1, the gain from this restricted reoptimization is exactly

$$\frac{1}{2} (\langle v, R_{-r}^*(x) \rangle - \lambda)_+^2.$$

530 The full K -slot optimum may further reoptimize the old coefficients, so the true improvement is at
531 least this large. Take expectations.

532 **E Ideal sparse-coding mechanism details**

533 The mechanism claims from Section 4 are made precise below: first-order merge pressure, event-
534 partition splitting, active-set maxima, and residual suppression. These are the exact formal statements
535 and proofs behind the mechanism narrative.

536 The same utility field creates merge pressure, split pressure, absorption-like residual suppression, and path-dependent alternatives.

537 **Deferred merge-pressure statement.** For $u \in S^{d-1}$ and $h \in u^\perp \cap S^{d-1}$, define the great-circle
 538 perturbation $v_\theta := \cos \theta u + \sin \theta h$.

Proposition E.1 (first-order merge pressure). *Let $u \in S^{d-1}$ and $h \in u^\perp \cap S^{d-1}$. Then*

$$539 \quad \left. \frac{d}{d\theta} \mathcal{U}_\lambda(v_\theta; R) \right|_{\theta=0} = \mathbb{E}[(\langle u, R \rangle - \lambda)_+ \langle h, R \rangle].$$

Consequently, if $\mathbb{E}[(\langle u, R \rangle - \lambda)_+ \langle h, R \rangle] > 0$, then u is not a local maximizer of $\mathcal{U}_\lambda(v; R)$ on S^{d-1} .

540 E.1 Proof of Proposition E.1

541 By Proposition D.2,

$$\nabla_v \mathcal{U}_\lambda(v; R) = \mathbb{E}[(\langle v, R \rangle - \lambda)_+ R].$$

542 Since v_θ is a great-circle perturbation through u with tangent h ,

$$\left. \frac{d}{d\theta} \mathcal{U}_\lambda(v_\theta; R) \right|_{\theta=0} = \langle \nabla_v \mathcal{U}_\lambda(u; R), h \rangle = \mathbb{E}[(\langle u, R \rangle - \lambda)_+ \langle h, R \rangle].$$

543 If the derivative is positive, utility increases along that tangent direction, so u is not a local maximizer.

544 **Deferred singleton-only statement.** The proposition below shows that merging is a mixed-event
 545 phenomenon, when higher-order support events are present.

Proposition E.2 (singleton-only populations do not strictly favor merges). *Let $x = AZ_1u + AZ_2c$, where $u \perp c$, $A > 0$, $Z_1, Z_2 \in \{0, 1\}$, and $Z_1Z_2 = 0$ a.s. Write $\pi_i := \mathbb{P}(Z_i = 1)$. If $0 < \lambda < A$, then*

$$\sup_{v \in S^{d-1}} \mathcal{U}_\lambda(v; x) = \frac{1}{2} \max\{\pi_1, \pi_2\} (A - \lambda)^2.$$

546 *A teacher direction with maximal singleton probability is globally optimal. If $\pi_1 = \pi_2$, both u and c are global maximizers, and no mixed direction strictly improves on them. If $\lambda = 0$, then*

$$\sup_{v \in S^{d-1}} \mathcal{U}_0(v; x) = \frac{1}{2} A^2 \max\{\pi_1, \pi_2\}.$$

When $\pi_1 = \pi_2$, every unit direction $v = au + bc$ with $a, b \geq 0$ and $a^2 + b^2 = 1$ is globally optimal, so mixed directions in the positive quadrant can tie but cannot strictly improve over teacher directions. Directions with a negative coordinate in the (u, c) basis are not optimal. If $\lambda \geq A$, then $\mathcal{U}_\lambda(v; x) = 0$ for every v .

547 E.2 Proof of Proposition E.2

548 For $v \in S^{d-1}$, write

$$a := \langle v, u \rangle, \quad b := \langle v, c \rangle, \quad \tau := \lambda/A.$$

549 Since $u \perp c$, we have

$$a^2 + b^2 \leq 1.$$

550 Because $Z_1Z_2 = 0$ almost surely,

$$\mathcal{U}_\lambda(v; x) = \frac{1}{2} A^2 \left[\pi_1 (a - \tau)_+^2 + \pi_2 (b - \tau)_+^2 \right].$$

551 Assume first that $0 < \tau < 1$. Set

$$y_1 := (a - \tau)_+, \quad y_2 := (b - \tau)_+.$$

552 If $y_1 = y_2 = 0$, the desired upper bound is immediate. Otherwise, on the active coordinates
 553 $I := \{i : y_i > 0\}$, we have $a_i = y_i + \tau$, where $a_1 = a$ and $a_2 = b$. Hence

$$\sum_{i \in I} (y_i + \tau)^2 \leq a^2 + b^2 \leq 1.$$

554 Therefore

$$\|y\|_2^2 + 2\tau \|y\|_1 + |I|\tau^2 \leq 1.$$

555 Since $|I| \geq 1$ and $\|y\|_1 \geq \|y\|_2$,

$$\|y\|_2^2 + 2\tau \|y\|_2 + \tau^2 \leq 1,$$

556 so

$$\|y\|_2 \leq 1 - \tau.$$

557 Thus

$$(a - \tau)_+^2 + (b - \tau)_+^2 = y_1^2 + y_2^2 \leq (1 - \tau)^2.$$

558 It follows that

$$\mathcal{U}_\lambda(v; x) \leq \frac{1}{2} A^2 \max\{\pi_1, \pi_2\} ((a - \tau)_+^2 + (b - \tau)_+^2) \leq \frac{1}{2} A^2 \max\{\pi_1, \pi_2\} (1 - \tau)^2.$$

559 Since $A^2(1 - \tau)^2 = (A - \lambda)^2$, this gives

$$\mathcal{U}_\lambda(v; x) \leq \frac{1}{2} \max\{\pi_1, \pi_2\} (A - \lambda)^2.$$

560 The bound is attained by $v = u$ if $\pi_1 = \max\{\pi_1, \pi_2\}$, and by $v = c$ if $\pi_2 = \max\{\pi_1, \pi_2\}$. This
 561 proves the claim for $0 < \lambda < A$.

562 If $\lambda = 0$, then

$$\mathcal{U}_0(v; x) = \frac{1}{2} A^2 [\pi_1 a_+^2 + \pi_2 b_+^2] \leq \frac{1}{2} A^2 \max\{\pi_1, \pi_2\} (a^2 + b^2) \leq \frac{1}{2} A^2 \max\{\pi_1, \pi_2\}.$$

563 The upper bound is attained by a highest-probability teacher direction. If $\pi_1 = \pi_2$, equality requires
 564 $a, b \geq 0$ and $a^2 + b^2 = 1$, which gives exactly the positive-quadrant unit directions in $\text{span}\{u, c\}$.

565 If $\lambda \geq A$, then for every unit v ,

$$A \langle v, u \rangle \leq A \leq \lambda, \quad A \langle v, c \rangle \leq A \leq \lambda,$$

566 so both positive parts vanish and $\mathcal{U}_\lambda(v; x) = 0$.

567 **Angle dependence of the merge mechanism.** The orthogonal construction isolates the threshold
 568 effect. If the same independent co-occurrence model is used with nonorthogonal unit directions u, c
 569 and

$$\rho := \langle u, c \rangle = \cos \theta \in (-1, 1),$$

570 then the normalized bisector is

$$m := \frac{u + c}{\|u + c\|} = \frac{u + c}{\sqrt{2(1 + \rho)}}.$$

571 On singleton events,

$$\langle m, Au \rangle = \langle m, Ac \rangle = A \sqrt{\frac{1 + \rho}{2}},$$

572 whereas on the joint event,

$$\langle m, A(u + c) \rangle = A \sqrt{2(1 + \rho)}.$$

573 Thus the threshold-only merge regime in which singleton teacher projections of size A are suppressed
 574 while the joint bisector projection remains active is

$$A \leq \lambda < A \sqrt{2(1 + \rho)}.$$

575 This interval is nonempty exactly when

$$\rho > -\frac{1}{2} \iff \theta < \frac{2\pi}{3}.$$

576 Hence, for fixed amplitude and threshold, more aligned features create stronger merge pressure, while
 577 sufficiently opposed features do not produce this threshold-only merge regime.

578 The local derivative also depends explicitly on the angle. Let

$$h := \frac{c - \rho u}{\sqrt{1 - \rho^2}},$$

579 the unit tangent at u pointing toward c , and define

$$v_\varphi := \cos \varphi u + \sin \varphi h.$$

580 For

$$x = As_1u + As_2c, \quad s_1, s_2 \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p),$$

581 the tangent derivative is

$$\left. \frac{d}{d\varphi} \mathcal{U}_\lambda(v_\varphi; x) \right|_{\varphi=0} = A\sqrt{1 - \rho^2} \left[p(1 - p)(A\rho - \lambda)_+ + p^2(A(1 + \rho) - \lambda)_+ \right].$$

582 For $\rho = 0$, this reduces to

$$Ap^2(A - \lambda)_+,$$

583 the orthogonal formula in [Proposition E.3](#).

584 **Deferred independent co-occurrence statement.** The proposition below shows that strict correla-
 585 tion is not necessary for merging to happen, since higher-order support events drive it.

Proposition E.3 (independent co-occurrence can favor a merge). *Let $x = As_1u + As_2c$, where
 $u \perp c$, $A > 0$, and $s_1, s_2 \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p)$. Let $m = (u + c)/\sqrt{2}$. If $A \leq \lambda < \sqrt{2}A$, then*

$$\mathcal{U}_\lambda(u; x) = 0 \quad \text{while} \quad \mathcal{U}_\lambda(m; x) > 0$$

586 *for every $p > 0$. If $\lambda < A$, then*

$$\left. \frac{d}{d\theta} \mathcal{U}_\lambda(\cos \theta u + \sin \theta c; x) \right|_{\theta=0} = Ap^2(A - \lambda) > 0.$$

587 E.3 Proof of [Proposition E.3](#)

588 The support events are $\{1\}$, $\{2\}$, and $\{1, 2\}$ with probabilities $p(1 - p)$, $p(1 - p)$, and p^2 . For u , only
 589 the events containing atom 1 matter, so

$$\mathcal{U}_\lambda(u; x) = \frac{1}{2}p(A - \lambda)_+^2.$$

590 For $m = (u + c)/\sqrt{2}$, singleton events contribute overlap $A/\sqrt{2}$ and the pair event contributes
 591 overlap $\sqrt{2}A$, yielding the stated formula by direct substitution into [Theorem 3.1](#). The derivative
 592 toward c follows from [Proposition E.1](#): only the joint event has nonzero tangent projection, giving

$$\left. \frac{d}{d\theta} \mathcal{U}_\lambda(\cos \theta u + \sin \theta c; x) \right|_{\theta=0} = Ap^2(A - \lambda)_+.$$

593 **Deferred absorption statements.** For an event A , write $\mathcal{U}_{\lambda,A}(v; R) := \frac{1}{2}\mathbb{E}[(\langle v, R \rangle - \lambda)_+^2 \mathbf{1}_A]$.

Theorem E.4 (post-learning residual suppression on a rank-one event). *Let A be an event
 on which $R = \xi u$, with $\|u\| = 1$ and $\xi \geq 0$. Let $v \in S^{d-1}$ be a previously selected feature
 with overlap $\rho := \langle u, v \rangle \in [0, 1]$ and define $\mathcal{R}_v(R) := R - v(\langle v, R \rangle - \lambda)_+$. If $\rho\xi > \lambda$ almost
 594 surely on A , then*

$$\mathcal{U}_{\lambda,A}(u; \mathcal{R}_v(R)) = \frac{1}{2}(1 - \rho)^2\mathbb{E}[(1 + \rho)\xi - \lambda]^2 \mathbf{1}_A.$$

Proposition E.5 (absorption-like false negatives from independent support intersections). *Let $u, w \in S^{d-1}$ satisfy $u \perp w$, and let $x = A(s_1 u + s_2 w)$, $s_1, s_2 \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p)$, with $A > 0$. Let $m = (u + w)/\sqrt{2}$, and assume $\frac{1}{\sqrt{2}} < \frac{\lambda}{A} < 1$. Then there exists $p_*(\lambda/A) < 1$ such that, for every $p > p_*(\lambda/A)$, m is the unique global maximizer of $v \mapsto \mathcal{U}_\lambda(v; x)$. After subtracting the oracle code along m , $R^+ := x - m(\langle m, x \rangle - \lambda)_+$, we have*

$$\mathcal{U}_\lambda(u; x) = \frac{1}{2}p(A - \lambda)^2, \quad \mathcal{U}_\lambda(u; R^+) = \frac{1}{2}p(1 - p)(A - \lambda)^2.$$

Eventwise, the later oracle gate for u fires on $\{s_1 = 1, s_2 = 0\}$ and not on $\{s_1 = s_2 = 1\}$. Hence, conditional on the teacher event $\{s_1 = 1\}$, its recall is $1 - p$.

596 E.4 Proofs for Section 4.4

597 *Proof of Theorem E.4.* On A , we have $R = \xi u$, and therefore

$$\langle v, R \rangle = \xi \langle v, u \rangle = \rho \xi.$$

598 By assumption, $\rho \xi > \lambda$ almost surely on A , so

$$(\langle v, R \rangle - \lambda)_+ = \rho \xi - \lambda.$$

599 Hence, on A ,

$$\mathcal{R}_v(R) = \xi u - v(\rho \xi - \lambda).$$

600 Taking the inner product with u gives

$$\begin{aligned} \langle u, \mathcal{R}_v(R) \rangle &= \xi - \rho(\rho \xi - \lambda) \\ &= (1 - \rho^2)\xi + \rho \lambda. \end{aligned}$$

601 Thus

$$\begin{aligned} \langle u, \mathcal{R}_v(R) \rangle - \lambda &= (1 - \rho^2)\xi - (1 - \rho)\lambda \\ &= (1 - \rho)((1 + \rho)\xi - \lambda). \end{aligned}$$

602 Since $\rho \xi > \lambda$ on A , we have $(1 + \rho)\xi - \lambda > 0$ on A . Taking the positive part and squaring yields

$$(\langle u, \mathcal{R}_v(R) \rangle - \lambda)_+^2 = (1 - \rho)^2((1 + \rho)\xi - \lambda)^2 \quad \text{on } A.$$

603 Multiplying by $\mathbf{1}_A$, taking expectation, and dividing by 2 proves the claim. \square

604 *Proof of Proposition E.5.* Write

$$t := \frac{\lambda}{A}, \quad a := \langle v, u \rangle, \quad b := \langle v, w \rangle.$$

605 Only the projection of v onto $\text{span}\{u, w\}$ affects the utility, so (a, b) ranges over the closed unit disk

$$D := \{(a, b) : a^2 + b^2 \leq 1\}.$$

606 For $v \in S^{d-1}$, the utility equals

$$\begin{aligned} \mathcal{U}_\lambda(v; x) &= \frac{1}{2}p(1 - p)(Aa - \lambda)_+^2 + \frac{1}{2}p(1 - p)(Ab - \lambda)_+^2 \\ &\quad + \frac{1}{2}p^2(A(a + b) - \lambda)_+^2. \end{aligned}$$

607 Equivalently, up to the positive factor $A^2/2$, define

$$F_p(a, b) := p(1 - p)(a - t)_+^2 + p(1 - p)(b - t)_+^2 + p^2(a + b - t)_+^2.$$

608 At $p = 1$,

$$F_1(a, b) = (a + b - t)_+^2.$$

609 Over D , the linear functional $a + b$ is uniquely maximized at

$$(a, b) = \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right),$$

610 corresponding to $v = m = (u + w)/\sqrt{2}$. Hence m is the unique maximizer of F_1 .

611 Because D is compact and $F_p \rightarrow F_1$ uniformly as $p \rightarrow 1$, strict uniqueness of the maximizer of F_1
 612 implies that, for every sufficiently small neighborhood N of $(1/\sqrt{2}, 1/\sqrt{2})$, all maximizers of F_p lie
 613 in N for p sufficiently close to 1. Since $t > 1/\sqrt{2}$, choose N small enough that

$$a < t, \quad b < t \quad \text{for all } (a, b) \in N.$$

614 On this neighborhood the singleton terms vanish, and

$$F_p(a, b) = p^2(a + b - t)_+^2.$$

615 This function is uniquely maximized over D at

$$(a, b) = \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right).$$

616 Therefore there exists $p_*(t) < 1$ such that, for every $p \in (p_*(t), 1]$, m is the unique global maximizer
 617 of $v \mapsto \mathcal{U}_\lambda(v; x)$.

618 It remains to compute the effect of subtracting the oracle code along m . Before subtraction,

$$\langle u, x \rangle = As_1,$$

619 so, since $A > \lambda$,

$$\mathcal{U}_\lambda(u; x) = \frac{1}{2}\mathbb{P}(s_1 = 1)(A - \lambda)^2 = \frac{1}{2}p(A - \lambda)^2.$$

620 Now define

$$R^+ := x - m(\langle m, x \rangle - \lambda)_+.$$

621 Because $t > 1/\sqrt{2}$,

$$\langle m, Au \rangle = \langle m, Aw \rangle = \frac{A}{\sqrt{2}} < \lambda,$$

622 so the m -code is inactive on the singleton events. Because $t < 1$,

$$\langle m, A(u + w) \rangle = A\sqrt{2} > \lambda,$$

623 so the m -code is active on the coactivation event $C = \{s_1 = s_2 = 1\}$.

624 Thus on $\{s_1 = 1, s_2 = 0\}$, the residual remains

$$R^+ = Au,$$

625 and the later u -gate fires with value $A - \lambda$. On C ,

$$\begin{aligned} R^+ &= A(u + w) - m(A\sqrt{2} - \lambda) \\ &= A(u + w) - \frac{u + w}{\sqrt{2}}(A\sqrt{2} - \lambda) \\ &= \frac{\lambda}{\sqrt{2}}(u + w). \end{aligned}$$

626 Therefore

$$\langle u, R^+ \rangle = \frac{\lambda}{\sqrt{2}} < \lambda \quad \text{on } C,$$

627 so the later u -gate is zero on the coactivation event. On all remaining events, either $s_1 = 0$ or the
 628 u -gate is also zero. Hence

$$\mathcal{U}_\lambda(u; R^+) = \frac{1}{2}\mathbb{P}(s_1 = 1, s_2 = 0)(A - \lambda)^2 = \frac{1}{2}p(1 - p)(A - \lambda)^2.$$

629 Finally, conditional on the teacher event $\{s_1 = 1\}$, the later u -gate fires exactly when $s_2 = 0$. Since
 630 s_1 and s_2 are independent,

$$\mathbb{P}(s_2 = 0 \mid s_1 = 1) = 1 - p.$$

631 This proves the recall statement. □

632 **Deferred residual-path statement.** The proposition below shows a mechanism analogous to seed
 633 instability in our framework.

Proposition E.6 (population residual-path non-identifiability). *Let u_1, u_2, u_3 be orthonormal and define $q_{ij} := (u_i + u_j)/\sqrt{2}$. Suppose the residual takes values $\sqrt{2}Aq_{12}, \sqrt{2}Aq_{13}, \sqrt{2}Aq_{23}$ with probabilities $\pi_{12}, \pi_{13}, \pi_{23} > 0$. If $\beta := \lambda/\sqrt{2}A \in (\sqrt{3}/2, 1)$, then the global maximizers of $v \mapsto \mathcal{U}_\lambda(v; R)$ are exactly the pair directions q_{ij} whose probabilities π_{ij} are maximal. In particular, tied probabilities give multiple inequivalent non-teacher population maximizers.*

635 E.5 Proof of Proposition E.6

636 Let

$$\mathcal{P} := \{12, 13, 23\}.$$

637 For $a \in \mathcal{P}$, write q_a for the corresponding pair direction and π_a for its probability. For any unit vector
 638 v , set

$$t_a := \langle v, q_a \rangle.$$

639 The pair directions satisfy

$$\langle q_a, q_b \rangle = \frac{1}{2} \quad (a \neq b).$$

640 We first show that at most one t_a can exceed β . If $t_a > \beta$ and $t_b > \beta$ for $a \neq b$, then

$$2\beta < t_a + t_b = \langle v, q_a + q_b \rangle \leq \|q_a + q_b\| = \sqrt{2 + 2\langle q_a, q_b \rangle} = \sqrt{3},$$

641 contradicting $\beta > \sqrt{3}/2$.

642 Therefore at most one event contributes positive utility for any v . The utility is

$$\mathcal{U}_\lambda(v; R) = \frac{1}{2} \sum_{a \in \mathcal{P}} \pi_a (\sqrt{2}A t_a - \lambda)_+^2 = A^2 \sum_{a \in \mathcal{P}} \pi_a (t_a - \beta)_+^2.$$

643 Since at most one term can be active and $t_a \leq 1$,

$$\mathcal{U}_\lambda(v; R) \leq A^2 (1 - \beta)^2 \max_{a \in \mathcal{P}} \pi_a.$$

644 This bound is attained by $v = q_a$ for every a with

$$\pi_a = \max_{b \in \mathcal{P}} \pi_b.$$

645 Indeed, for $v = q_a$,

$$t_a = 1, \quad t_b = \frac{1}{2} < \beta \quad (b \neq a),$$

646 so only the a -event contributes. Hence the global maximizers are exactly the most probable pair
 647 directions.

648 The utility gap between two pair directions is explicit:

$$\mathcal{U}_\lambda(q_a; R) - \mathcal{U}_\lambda(q_b; R) = A^2 (1 - \beta)^2 (\pi_a - \pi_b).$$

649 Thus exact ties give exact population non-identifiability, and near ties give proportionally small
 650 population gaps.

651 Teacher directions have zero utility. For example,

$$\langle u_1, \sqrt{2}Aq_{12} \rangle = A, \quad \langle u_1, \sqrt{2}Aq_{13} \rangle = A, \quad \langle u_1, \sqrt{2}Aq_{23} \rangle = 0.$$

652 Since

$$\lambda > \frac{\sqrt{3}}{2} \sqrt{2}A = \sqrt{\frac{3}{2}}A > A,$$

653 all projections are below threshold. Thus

$$\mathcal{U}_\lambda(u_1; R) = 0,$$

654 and similarly for u_2, u_3 .

655 Now suppose q_a is selected first and the oracle code

$$z_a(R) := (\langle q_a, R \rangle - \lambda)_+$$

656 is subtracted. On the a -event,

$$R = \sqrt{2}Aq_a, \quad z_a = \sqrt{2}A - \lambda,$$

657 so the new residual is

$$R^+ = \lambda q_a.$$

658 For every unit direction v ,

$$\langle v, R^+ \rangle - \lambda = \lambda \langle v, q_a \rangle - \lambda \leq 0.$$

659 Thus the selected event contributes no future above-threshold utility to any direction. On every other
660 pair event $b \neq a$, the code is inactive because

$$\langle q_a, \sqrt{2}Aq_b \rangle = \frac{\sqrt{2}A}{2} < \lambda,$$

661 where the last inequality follows from $\beta > \sqrt{3}/2 > 1/2$. Hence all unselected pair events are
662 unchanged. The later population utility field is therefore the original pair-event utility with the
663 selected event removed.

664 At $\beta = \sqrt{3}/2$, the strict exclusion argument becomes non-strict: two pair directions can reach the
665 threshold simultaneously, so the winner-takes-one proof no longer applies. Below this threshold,
666 simultaneous active pair events are possible.

667 **Deferred splitting statement.** The proposition below shows a role of event partition in splitting.

Proposition E.7 (event-partition splitting certificate). *Suppose an event family A decomposes into m disjoint rank-one subevents A_1, \dots, A_m , with $R_B = \xi_k u_k$ on A_k , where each $u_k \in S^{d-1}$ and each ξ_k is a nonnegative random scalar. If the best one-direction utility on A is strictly smaller than the utility achieved by the specialized directions u_1, \dots, u_m , then the oracle objective strictly favors allocating separate directions to the subevents.*

669 E.6 Proof of Proposition E.7

670 The exact sufficient condition used by Proposition E.7 is

$$\sup_{v \in S^{d-1}} \sum_{k=1}^m \mathbb{E}[(\xi_k \langle v, u_k \rangle - \lambda)_+^2 \mathbf{1}_{A_k}] < \sum_{k=1}^m \mathbb{E}[(\xi_k - \lambda)_+^2 \mathbf{1}_{A_k}],$$

671 where $A = \bigsqcup_{k=1}^m A_k$, $R_B = \xi_k u_k$ a.s. on A_k , $\xi_k \geq 0$, and $\|u_k\| = 1$. This is the formal version of
672 event-partition splitting: one broad direction is worse than specialized children when no single v can
673 match the eventwise rank-one gains.

674 On each subevent A_k , the residual is $\xi_k u_k$. A single direction v therefore achieves eventwise gain

$$\frac{1}{2} \sum_{k=1}^m \mathbb{E}[(\xi_k \langle v, u_k \rangle - \lambda)_+^2 \mathbf{1}_{A_k}].$$

675 By contrast, the child family $\{u_k\}$ achieves gain at least

$$\frac{1}{2} \sum_{k=1}^m \mathbb{E}[(\xi_k - \lambda)_+^2 \mathbf{1}_{A_k}]$$

676 by the feasible code $z_k = (\xi_k - \lambda)_+ \mathbf{1}_{A_k}$. If the former is strictly smaller than the latter uniformly
677 over v , then splitting is oracle-preferred.

678 **Deferred two-subevent splitting statement.** Proposition below nuances the viewpoint above.

Proposition E.8 (two-subevent splitting regime). *Let $A = A_1 \sqcup A_2$, with $\mathbb{P}(A_1) = \mathbb{P}(A_2) = q/2$. Suppose $R_B = A_0 u_k$ on A_k , where $A_0 > 0$, $u_1, u_2 \in S^{d-1}$, and $\angle(u_1, u_2) = \theta > 0$. If $0 \leq \lambda < A_0$, then the two specialized directions u_1, u_2 achieve strictly larger eventwise oracle utility on A than any single shared direction. With $\tau = \lambda/A_0$, their guaranteed eventwise advantage is at least*

679

$$\frac{q}{2}(A_0 - \lambda)^2 - \sup_{v \in S^{d-1}} \frac{q}{4} \sum_{k=1}^2 (A_0 \langle v, u_k \rangle - \lambda)_+^2 \geq \frac{qA_0^2}{4} [(1 - \tau)^2 - (\cos(\theta/2) - \tau)_+]^2 > 0.$$

680 **E.7 Proof of Proposition E.8**

681 The specialized directions u_1, u_2 achieve at least the eventwise utility

$$\frac{1}{2} \sum_{k=1}^2 \mathbb{E}[(A_0 - \lambda)^2 \mathbf{1}_{A_k}] = \frac{q}{2}(A_0 - \lambda)^2.$$

682 For a single shared direction v , write $t_k = \langle v, u_k \rangle$ and $\tau = \lambda/A_0$. The eventwise utility is

$$\frac{qA_0^2}{4} [(t_1 - \tau)_+^2 + (t_2 - \tau)_+^2].$$

683 For every unit v , at least one of t_1, t_2 is at most $\cos(\theta/2)$. Otherwise

$$\langle v, u_1 + u_2 \rangle > 2 \cos(\theta/2) = \|u_1 + u_2\|,$$

684 which contradicts Cauchy–Schwarz. Since also $t_k \leq 1$, we have

$$(t_1 - \tau)_+^2 + (t_2 - \tau)_+^2 \leq (1 - \tau)^2 + (\cos(\theta/2) - \tau)_+^2.$$

685 Subtracting this upper bound from the specialized utility gives the displayed margin. The margin is
686 strictly positive because $0 \leq \tau < 1$ and $\cos(\theta/2) < 1$.

687 For the optional active-set characterization below, specialize to an atomic residual law

$$R = Au_i \quad \text{with probability } p_i, \quad i = 1, \dots, M,$$

688 where $A > 0$, $u_i \in S^{d-1}$, $p_i > 0$, and $\sum_i p_i = 1$. Write $\beta := \lambda/A$. For $S \subseteq [M]$, define the active
689 cell

$$\mathcal{C}_S := \{v \in S^{d-1} : \langle v, u_i \rangle > \beta \text{ if and only if } i \in S\},$$

690 and the weighted active-set moments

$$G_S := \sum_{i \in S} p_i u_i u_i^\top, \quad b_S := \sum_{i \in S} p_i u_i.$$

Theorem E.9 (exact active-set characterization of nonboundary strict local maxima). Assume $d \geq 2$ and the atomic residual law

$$\mathbb{P}(R = Au_i) = p_i, \quad i = 1, \dots, M,$$

where $A > 0$, $u_i \in S^{d-1}$, $p_i > 0$, and $\sum_{i=1}^M p_i = 1$. Set

$$\beta := \frac{\lambda}{A}.$$

For $S \subseteq [M]$, define the strict active cell

$$\mathcal{C}_S^\circ := \{v \in S^{d-1} : \langle v, u_i \rangle > \beta \forall i \in S, \quad \langle v, u_i \rangle < \beta \forall i \notin S\}.$$

For $S \neq \emptyset$, define

$$G_S := \sum_{i \in S} p_i u_i u_i^\top, \quad b_S := \sum_{i \in S} p_i u_i.$$

691 Let $S \neq \emptyset$ and $v \in \mathcal{C}_S^\circ$. Then, on a neighborhood of v in S^{d-1} , the active set is fixed and

$$\mathcal{U}_\lambda(v; R) = \frac{A^2}{2} \left(v^\top G_S v - 2\beta b_S^\top v + \beta^2 \sum_{i \in S} p_i \right).$$

Moreover, v is a spherical critical point if and only if there exists $\eta \in \mathbb{R}$ such that

$$G_S v - \beta b_S = \eta v.$$

At such a critical point, v is a strict local maximizer of \mathcal{U}_λ on S^{d-1} if and only if

$$\xi^\top G_S \xi < \eta \|\xi\|^2 \quad \forall \xi \in v^\perp \setminus \{0\}.$$

Equivalently,

$$\eta > \lambda_{\max}(P_{v^\perp} G_S P_{v^\perp} |_{v^\perp}).$$

Thus every nonboundary strict local maximizer is a threshold-consistent stable solution of the shifted active-set eigenproblem

$$G_S v - \beta b_S = \eta v.$$

692 E.8 Proof of Theorem E.9

693 Since $v \in \mathcal{C}_S^\circ$, the active inequalities are strict. Hence, on a neighborhood of v in S^{d-1} , the active set
694 is fixed and

$$\mathcal{U}_\lambda(v; R) = \frac{A^2}{2} \sum_{i \in S} p_i (\langle v, u_i \rangle - \beta)^2.$$

695 Expanding gives

$$\mathcal{U}_\lambda(v; R) = \frac{A^2}{2} \left(v^\top G_S v - 2\beta b_S^\top v + \beta^2 \sum_{i \in S} p_i \right).$$

696 Thus, inside the active cell,

$$\nabla_v \mathcal{U}_\lambda(v; R) = A^2 (G_S v - \beta b_S).$$

697 The spherical critical-point condition is

$$P_{v^\perp} \nabla_v \mathcal{U}_\lambda(v; R) = 0,$$

698 which is equivalent to the existence of $\eta \in \mathbb{R}$ such that

$$G_S v - \beta b_S = \eta v.$$

699 At such a critical point, the Euclidean Hessian inside the active cell is $A^2 G_S$. Therefore the
700 Riemannian Hessian on S^{d-1} , restricted to v^\perp , is

$$\text{Hess}_{S^{d-1}} \mathcal{U}_\lambda(v) = A^2 (P_{v^\perp} G_S P_{v^\perp} - \eta I_{v^\perp}).$$

701 Hence the Riemannian Hessian is negative definite if and only if

$$\xi^\top G_S \xi < \eta \|\xi\|^2 \quad \forall \xi \in v^\perp \setminus \{0\}.$$

702 Negative definiteness is sufficient for v to be a strict local maximizer.

703 It remains to show necessity. Suppose v is a strict local maximizer. Then the Riemannian Hessian is
704 negative semidefinite, so

$$\xi^\top G_S \xi \leq \eta \|\xi\|^2 \quad \forall \xi \in v^\perp.$$

705 Assume, for contradiction, that equality holds for some $\xi \in v^\perp$ with $\|\xi\| = 1$. Consider the spherical
706 geodesic

$$v(t) = v \cos t + \xi \sin t.$$

707 For sufficiently small $|t|$, $v(t) \in \mathcal{C}_S^\circ$. Define the rescaled active-cell objective, up to an irrelevant
708 additive constant, by

$$f(w) := w^\top G_S w - 2\beta b_S^\top w.$$

709 Using $G_S v - \beta b_S = \eta v$, $\xi \perp v$, and $\xi^\top G_S \xi = \eta$, a direct expansion gives

$$f(v(t)) - f(v) = \beta(b_S^\top v)(1 - \cos t)^2 - 2\beta(b_S^\top \xi) \sin t(1 - \cos t).$$

710 If $\beta b_S^\top \xi \neq 0$, the second term has leading order

$$-\beta(b_S^\top \xi)t^3,$$

711 which changes sign with t . Hence v cannot be a local maximizer.

712 If $\beta b_S^\top \xi = 0$, then either $\beta = 0$ or $b_S^\top \xi = 0$. If $\beta = 0$, the expression above is identically zero along
713 this geodesic, so v cannot be a strict local maximizer. If $\beta > 0$, then $b_S^\top \xi = 0$, and

$$f(v(t)) - f(v) = \beta(b_S^\top v)(1 - \cos t)^2.$$

714 Since $v \in \mathcal{C}_S^\circ$,

$$b_S^\top v = \sum_{i \in S} p_i \langle v, u_i \rangle > \beta \sum_{i \in S} p_i > 0.$$

715 Therefore $f(v(t)) - f(v) > 0$ for all sufficiently small $t \neq 0$, again contradicting local maximality.

716 Thus equality cannot occur in any nonzero tangent direction. The Hessian is negative definite, proving
717 the strict inequality and completing the characterization.

718 E.9 Singleton-cell and pairwise-exclusion lemmas

719 If $\max_{j \neq i} \langle u_i, u_j \rangle < \beta$, then at u_i one has $\langle u_i, u_i \rangle = 1 > \beta$ and $\langle u_i, u_j \rangle < \beta$ for every $j \neq i$. These
720 inequalities persist in a neighborhood, so u_i lies in the interior of the singleton cell.

721 For pairwise exclusion, suppose two atoms u_i, u_j were simultaneously active at some v , so that
722 $\langle v, u_i \rangle > \beta$ and $\langle v, u_j \rangle > \beta$. Then

$$\langle v, u_i + u_j \rangle > 2\beta.$$

723 But

$$\langle v, u_i + u_j \rangle \leq \|u_i + u_j\| = \sqrt{2 + 2\langle u_i, u_j \rangle} \leq 2\nu,$$

724 contradicting $\beta > \nu$, where

$$\nu := \max_{i \neq j} \sqrt{\frac{1 + \langle u_i, u_j \rangle}{2}}.$$

725 E.10 Split rank as an architectural refinement

726 The oracle mechanisms above interact with an architectural mechanism specific to amortized SAEs:
727 the oracle gate associated with a decoder direction may fail to lie in the single affine-ReLU class.
728 This motivates a simple complexity notion.

Definition E.10 (split rank). For a nonnegative measurable target $g : \mathbb{R}^d \rightarrow \mathbb{R}_+$, define

$$\text{srank}(g) := \min \left\{ m \in \mathbb{N} : g(\nu) = \sum_{k=1}^m (e_k^\top \nu + b_k)_+ \text{ a.s.} \right\},$$

when such an m exists, and set $\text{srank}(g) = \infty$ otherwise.

If the oracle gate of a decoder direction has split rank greater than one, then a standard one-ReLU feature must either incur amortization error or factor the target across several slots. This mechanism is distinct from oracle event-partition splitting, which concerns which support partition is reward-optimal even with unrestricted activations.

The split-rank definition is therefore an architectural refinement, not another primitive oracle mechanism. It explains why the same oracle direction can appear as several learned slots when a single affine-ReLU gate cannot express the event rule.

F Positive and robust recovery refinements

The counterexamples are complemented by sufficient conditions: singleton mass can dominate contamination, directional moments can remain stable under small support errors, and low coherence can control off-target signal. These results delimit the regimes in which the usual feature-recovery intuition is scientifically justified.

Singleton dominance can survive margins, perturbations, low coherence, and random or learned geometry assumptions.

F.1 Proof of Theorem 4.1

The quantities used in the theorem are

$$\Delta_{j,\lambda}^B(\varepsilon) := \frac{1}{2} \mathbb{E} \left[\left((\alpha_j \|\tilde{w}_j\| - \lambda)_+^2 - (\alpha_j \|\tilde{w}_j\| \cos \varepsilon - \lambda)_+^2 \right) \mathbf{1}_{E_{\{j\}}(B)} \right],$$

and

$$C_j^B := \frac{1}{2} \mathbb{E} [\|R_B\|^2 \mathbf{1}_{E_{\{j\}}(B)^c}].$$

On $E_{\{j\}}(B)$,

$$R_B = \alpha_j \tilde{w}_j,$$

so for every v with $\angle(v, u_j) \geq \varepsilon$,

$$\langle u_j, R_B \rangle = \alpha_j \|\tilde{w}_j\|, \quad \langle v, R_B \rangle \leq \alpha_j \|\tilde{w}_j\| \cos \varepsilon.$$

Therefore the utility difference restricted to $E_{\{j\}}(B)$ is at least $\Delta_{j,\lambda}^B(\varepsilon)$. On the complement $E_{\{j\}}(B)^c$, competitor utility is bounded by

$$\frac{1}{2} \mathbb{E} [\|R_B\|^2 \mathbf{1}_{E_{\{j\}}(B)^c}] = C_j^B,$$

while the utility of u_j is nonnegative. Subtracting the two bounds yields the claim.

Sharper threshold-aware competitor load. The proof above uses the threshold-blind off-singleton load C_j^B . A sharper but less explicit quantity is

$$C_{j,\lambda}^B(\varepsilon) := \sup_{\angle(v, u_j) \geq \varepsilon} \frac{1}{2} \mathbb{E} \left[\left(\langle v, R_B \rangle - \lambda \right)_+^2 \mathbf{1}_{E_{\{j\}}(B)^c} \right].$$

The same argument proves localization under the weaker condition

$$\Delta_{j,\lambda}^B(\varepsilon) > C_{j,\lambda}^B(\varepsilon).$$

754 Indeed, for every v with $\angle(v, u_j) \geq \varepsilon$, the singleton-event contribution to

$$\mathcal{U}_\lambda(u_j; R_B) - \mathcal{U}_\lambda(v; R_B)$$

755 is at least $\Delta_{j,\lambda}^B(\varepsilon)$, while the off-singleton contribution of v is at most $C_{j,\lambda}^B(\varepsilon)$, and the off-singleton
756 contribution of u_j is nonnegative. Hence

$$\mathcal{U}_\lambda(u_j; R_B) - \mathcal{U}_\lambda(v; R_B) \geq \Delta_{j,\lambda}^B(\varepsilon) - C_{j,\lambda}^B(\varepsilon) > 0.$$

757 Thus no maximizer can lie outside the ε -ball around u_j .

758 The simpler load C_j^B used in the main theorem is a conservative upper bound on $C_{j,\lambda}^B(\varepsilon)$. For every
759 unit v ,

$$\langle v, R_B \rangle - \lambda \leq (\langle v, R_B \rangle)_+^2 \leq \langle v, R_B \rangle^2 \leq \|R_B\|^2,$$

760 and therefore

$$C_{j,\lambda}^B(\varepsilon) \leq \frac{1}{2} \mathbb{E}[\|R_B\|^2 \mathbf{1}_{E_{\{j\}}(B)^c}] = C_j^B.$$

761 Thus the main-text condition

$$\Delta_{j,\lambda}^B(\varepsilon) > C_j^B$$

762 is a simple sufficient condition for the sharper threshold-aware condition.

Corollary F.1 (uniform singleton-margin criterion). *Suppose that on $E_{\{j\}}(B)$,*

$$\alpha_j \|\tilde{w}_j\| \geq a_j > \lambda \quad \text{a.s.},$$

and write

$$q_j^B := \mathbb{P}(E_{\{j\}}(B)).$$

Then

763
$$\Delta_{j,\lambda}^B(\varepsilon) \geq \frac{q_j^B}{2} \left[(a_j - \lambda)^2 - (a_j \cos \varepsilon - \lambda)_+^2 \right].$$

Hence the sufficient condition

$$q_j^B \left[(a_j - \lambda)^2 - (a_j \cos \varepsilon - \lambda)_+^2 \right] > \mathbb{E}[\|R_B\|^2 \mathbf{1}_{E_{\{j\}}(B)^c}]$$

localizes all global maximizers near u_j .

764 F.2 Proof of Corollary F.1

765 On $E_{\{j\}}(B)$, the quantity $\alpha_j \|\tilde{w}_j\|$ is bounded below by $a_j > \lambda$ almost surely. Therefore the
766 integrand defining $\Delta_{j,\lambda}^B(\varepsilon)$ is bounded below by

$$(a_j - \lambda)^2 - (a_j \cos \varepsilon - \lambda)_+^2.$$

767 Multiplying by the event probability q_j^B and the prefactor 1/2 gives the stated lower bound. The
768 localization condition then follows by combining this lower bound with [Theorem 4.1](#).

Proposition F.2 (robust witness-concentration moment bound). *Let A be a measurable event
and suppose that for some unit u and activation $a \geq 0$ with $\mathbb{E}[\|R_B\| | a] < \infty$,*

$$\mu(a; B) = \mathbb{E}[R_B a] = \beta_A(a)u + \eta_A(a), \quad \eta_A(a) \in u^\perp, \quad \beta_A(a) > 0.$$

Then

769
$$\tan \angle(\mu(a; B), u) = \frac{\|\eta_A(a)\|}{\beta_A(a)}.$$

If $\beta_A(a) > 2 \|\eta_A(a)\|$, then

$$\left\| \frac{\mu(a; B)}{\|\mu(a; B)\|} - u \right\| \leq \frac{2 \|\eta_A(a)\|}{\beta_A(a)}.$$

770 **F.3 Proof of Proposition F.2**

771 Write

$$\mu(a; B) = \beta_A(a)u + \eta_A(a).$$

772 Since $\eta_A(a) \in u^\perp$,

$$\|P_{u^\perp}\mu(a; B)\| = \|\eta_A(a)\|, \quad \langle u, \mu(a; B) \rangle = \beta_A(a).$$

773 Taking the ratio of the transverse norm and the longitudinal component gives the tangent identity.
 774 The normalized-vector bound follows from the assumption $\beta_A(a) > 2\|\eta_A(a)\|$ and a standard
 775 perturbation estimate.

776 **F.4 Support is non-robust while direction is robust**

Proposition F.3 (support is non-robust while direction is robust). *Let A be a rank-one event for the ray \mathbb{R}_+u , and let $a \geq 0$, with $\mathbb{E}[\|R_B\|a] < \infty$, satisfy*

$$a = 0 \quad \text{a.s. on } A^c, \quad \mu(a; B) \neq 0.$$

Let $h \geq 0$ be any measurable function supported on A^c and nonzero on a positive-probability subset of A^c , with $\mathbb{E}[\|R_B\|h] < \infty$. Define

$$a_\varepsilon := a + \varepsilon h, \quad \varepsilon > 0.$$

777 *Then for every $\varepsilon > 0$, the activation a_ε no longer fires exactly on A . However,*

$$\mu(a_\varepsilon; B) = \mu(a; B) + \varepsilon \mathbb{E}[R_B h],$$

and if

$$\varepsilon \mathbb{E}[\|R_B\|h] \leq \frac{1}{2} \|\mu(a; B)\|,$$

then

$$\left\| \frac{\mu(a_\varepsilon; B)}{\|\mu(a_\varepsilon; B)\|} - \frac{\mu(a; B)}{\|\mu(a; B)\|} \right\| \leq \frac{4\varepsilon \mathbb{E}[\|R_B\|h]}{\|\mu(a; B)\|}.$$

778 *Proof.* Exact support recovery fails immediately because h is supported on A^c and is nonzero there.
 779 The moment identity is linear in the activation. The perturbation satisfies

$$\|\varepsilon \mathbb{E}[R_B h]\| \leq \varepsilon \mathbb{E}[\|R_B\|h].$$

780 Under the assumed bound, this perturbation has norm at most $\|\mu(a; B)\|/2$, so the standard
 781 normalized-vector perturbation inequality yields the stated estimate. \square

782 **E.5 Low coherence controls off-target contamination**

Proposition F.4 (low coherence controls off-target contamination). *Let $w_1^*, \dots, w_F^* \in S^{d-1}$ be teacher directions with coherence*

$$\gamma := \max_{i \neq j} |\langle w_i^*, w_j^* \rangle|.$$

Assume

783
$$x = \sum_{i=1}^F f_i w_i^*, \quad f_i \geq 0, \quad |\{i : f_i > 0\}| \leq s, \quad 0 \leq f_i \leq a_{\max}.$$

Then for every j ,

$$\left| \sum_{i \neq j} \langle w_j^*, w_i^* \rangle f_i \right| \leq \gamma(s-1)a_{\max},$$

and therefore

$$\langle w_j^*, x \rangle \geq f_j - \gamma(s-1)a_{\max}.$$

784 *Proof.* For every j ,

$$\left| \sum_{i \neq j} \langle w_j^*, w_i^* \rangle f_i \right| \leq \sum_{i \neq j} |\langle w_j^*, w_i^* \rangle| f_i \leq \gamma \sum_{i \neq j} f_i \leq \gamma(s-1)a_{\max}.$$

785 Since

$$\langle w_j^*, x \rangle = f_j + \sum_{i \neq j} \langle w_j^*, w_i^* \rangle f_i,$$

786 the lower bound follows. □

Proposition E.5 (random learned geometry is low-coherence with high probability). *Let w_1^*, \dots, w_F^* be i.i.d. uniform on S^{d-1} . Then there exists a universal constant $c > 0$ such that, with probability at least $1 - \delta$,*

787

$$\max_{i \neq j} |\langle w_i^*, w_j^* \rangle| \leq c \sqrt{\frac{\log(F/\delta)}{d}}.$$

788 *Proof.* For a fixed pair (i, j) , rotational invariance and spherical concentration imply

$$\mathbb{P}(|\langle w_i^*, w_j^* \rangle| > t) \leq 2e^{-c_0 dt^2}$$

789 for a universal constant $c_0 > 0$ and all $0 \leq t \leq 1$. Choose

$$t = C \sqrt{\frac{\log(F/\delta)}{d}}$$

790 with C large enough that $2 \binom{F}{2} e^{-c_0 dt^2} \leq \delta$. A union bound over the $\binom{F}{2}$ pairs yields

$$\max_{i \neq j} |\langle w_i^*, w_j^* \rangle| \leq C \sqrt{\frac{\log(F/\delta)}{d}}$$

791 with probability at least $1 - \delta$. Absorbing C into the universal constant c gives the claim. □

792 This is a standard concentration fact included only to make the low-coherence positive regime explicit.

793 **G Learned teacher geometry**

794 If the teacher dictionary is itself learned by an upstream sparse-superposition objective, the oracle
 795 picture gains one useful nuance: upstream geometry can make support-to-direction alignment easier
 796 in special symmetric regimes, but exact reconstruction of rich support families forces strong Gram-
 797 structure constraints.

798 Symmetry and simplex-like learned geometry can cancel off-target moments, so support-faithful gates
 can recover directions even outside literal singleton-only settings.

799 Let

$$f \in \mathbb{R}_+^F$$

800 be drawn from a sparse generator, and consider the tied feature autoencoder

$$\hat{f} = \sigma(Gf + b), \quad \sigma(t) = t_+, \quad G = WW^\top \succeq 0, \quad \text{rank}(G) \leq d,$$

801 trained with

$$\mathcal{J}(G, b) = \mathbb{E}[\|f - \sigma(Gf + b)\|_2^2].$$

802 The learned teacher directions are the rows of W , denoted

$$w_1^*, \dots, w_F^* \in \mathbb{R}^d,$$

803 so that

$$G_{ij} = \langle w_i^*, w_j^* \rangle.$$

804 Downstream, the hidden state is

$$x = \sum_{i=1}^F f_i w_i^*.$$

805 For a support family $T \subseteq [F]$, write

$$E_T := \{\text{supp}(f) = T\},$$

806 and let $\Omega_T \subseteq (0, \infty)^T$ be the topological support of f_T conditional on E_T . We say that E_T is
 807 full-dimensional if Ω_T contains a nonempty open subset of \mathbb{R}^T .

Theorem G.1 (active-block rigidity of the learned teacher). *Assume E_T is full-dimensional and the feature autoencoder reconstructs it exactly:*

$$\sigma(Gf + b) = f \quad \text{a.s. on } E_T.$$

Then

$$G_{TT} = I, \quad b_T = 0.$$

808 *Moreover, for every $j \notin T$,*

$$G_{jTs} + b_j \leq 0 \quad \forall s \in \Omega_T.$$

Consequently, if every pair support $E_{\{i,j\}}$ is full-dimensional and reconstructed exactly, then

$$G_{ij} = 0 \quad \forall i \neq j.$$

If this holds for all pairs and $\text{rank}(G) \leq d < F$, exact reconstruction is impossible.

809 *Proof.* On the support event E_T , one has $f_{T^c} = 0$. Exact reconstruction gives

$$\sigma(G_{TT}f_T + b_T) = f_T, \quad \sigma(G_{T^cT}f_T + b_{T^c}) = 0$$

810 almost surely conditional on E_T .

811 Because $f_T \in (0, \infty)^T$ on E_T , the first identity implies

$$G_{TT}f_T + b_T = f_T$$

812 almost surely on E_T . The map

$$s \mapsto (G_{TT} - I)s + b_T$$

813 is continuous. Since the conditional topological support Ω_T contains a nonempty open set, and the
 814 continuous map vanishes almost surely, it vanishes on Ω_T , hence on a nonempty open set. An affine
 815 map that vanishes on a nonempty open set is identically zero. Therefore

$$G_{TT} = I, \quad b_T = 0.$$

816 For $j \notin T$, exact inactivity gives

$$G_{jT}f_T + b_j \leq 0$$

817 almost surely on E_T . By the same continuity and support argument,

$$G_{jT}s + b_j \leq 0 \quad \forall s \in \Omega_T.$$

818 If every pair event $E_{\{i,j\}}$ is full-dimensional and reconstructed exactly, applying the first part with
 819 $T = \{i, j\}$ gives

$$G_{\{i,j\},\{i,j\}} = I_2,$$

820 so

$$G_{ij} = G_{ji} = 0.$$

821 If this holds for all pairs, then $G = I_F$. But $\text{rank}(I_F) = F$, contradicting $\text{rank}(G) \leq d < F$. Hence
 822 exact reconstruction is impossible in the undercomplete regime. \square

Lemma G.2 (singleton sign constraint). *Assume each singleton event $E_{\{i\}}$ has at least two positive amplitudes in its conditional support and is reconstructed exactly. Then*

$$G_{ii} = 1, \quad b_i = 0, \quad G_{ji} \leq 0 \quad (j \neq i).$$

824 *Proof of Lemma G.2.* Assume each singleton event $E_{\{i\}}$ has at least two positive amplitudes in
 825 its conditional support and is reconstructed exactly. On $E_{\{i\}}$, exact reconstruction of the active
 826 coordinate gives

$$(G_{ii}f_i + b_i)_+ = f_i.$$

827 Since $f_i > 0$, this implies

$$G_{ii}f_i + b_i = f_i$$

828 on at least two distinct positive values of f_i . Hence

$$G_{ii} = 1, \quad b_i = 0.$$

829 Now fix $j \neq i$. Exact inactivity of coordinate j on $E_{\{i\}}$ gives

$$G_{ji}f_i + b_j \leq 0.$$

830 Since the previous argument applied to singleton j gives $b_j = 0$, and since $f_i > 0$, we obtain

$$G_{ji} \leq 0.$$

831 \square

Theorem G.3 (permutation symmetry induces equiangular Gram blocks). *Let $T \subseteq [F]$, $|T| = k$. Assume the law of f is invariant under all permutations of the coordinates in T , and let (G, b) be a teacher solution with the same symmetry:*

$$PGP^\top = G, \quad Pb = b \quad \forall P \in \mathfrak{S}_T.$$

Then

$$G_{TT} = \alpha I + \beta \mathbf{1}\mathbf{1}^\top, \quad b_T = \tau \mathbf{1}$$

for some scalars α, β, τ .

832

If, in addition, all singleton events $E_{\{i\}}$, $i \in T$, are nondegenerate and reconstructed exactly, then

$$b_T = 0, \quad G_{TT} = (1 - \rho)I + \rho \mathbf{1}\mathbf{1}^\top$$

for some $\rho \leq 0$. If also $\text{rank}(G_{TT}) < k$, then

$$\rho = -\frac{1}{k-1},$$

so G_{TT} is the Gram matrix of a regular simplex.

833 *Proof.* A matrix commuting with every permutation matrix on T lies in the commutant of the
834 permutation representation of \mathfrak{S}_T . This commutant consists exactly of matrices of the form

$$\alpha I + \beta \mathbf{1}\mathbf{1}^\top.$$

835 Similarly, the only vectors fixed by all permutations of T are constant vectors, so

$$b_T = \tau \mathbf{1}.$$

836 If all singleton events $E_{\{i\}}$, $i \in T$, are nondegenerate and reconstructed exactly, [Lemma G.2](#) gives

$$G_{ii} = 1, \quad b_i = 0, \quad G_{ij} \leq 0 \quad (i \neq j).$$

837 By symmetry, all off-diagonal entries in G_{TT} are equal. Denote their common value by $\rho \leq 0$. Then

$$G_{TT} = (1 - \rho)I + \rho \mathbf{1}\mathbf{1}^\top, \quad b_T = 0.$$

838 The eigenvalues of this matrix are

$$1 - \rho \quad \text{with multiplicity } k - 1,$$

839 and

$$1 + (k - 1)\rho \quad \text{with multiplicity } 1.$$

840 Since $G_{TT} \succeq 0$, these eigenvalues are nonnegative. Because $\rho \leq 0$, the eigenvalue $1 - \rho$ is strictly
841 positive. Therefore $\text{rank}(G_{TT}) < k$ can only occur when

$$1 + (k - 1)\rho = 0,$$

842 i.e.

$$\rho = -\frac{1}{k-1}.$$

843 This is the regular-simplex Gram matrix. □

Proposition G.4 (effective threshold in an equiangular learned block). *Let $T \subseteq [F]$ be a cluster of size k satisfying*

$$\|w_i^*\| = 1, \quad \langle w_i^*, w_j^* \rangle = \rho \quad (i \neq j, i, j \in T),$$

with

$$-\frac{1}{k-1} < \rho < 1.$$

844 Suppose the oracle nonnegative Lasso with fixed decoder W^* has active support exactly T , and there is no active off-cluster contribution into T , meaning $(W_T^*)^\top x = G_{TT}^* f_T$ on the event under consideration. Then, for every $i \in T$,

$$z_i^* = f_i - \frac{\lambda}{1 + (k-1)\rho}.$$

Thus the effective threshold inside the block is

$$\lambda_{\text{eff}}(k, \rho) = \frac{\lambda}{1 + (k-1)\rho}.$$

845 *Proof.* On active support T , the KKT equations for the nonnegative Lasso are

$$G_{TT}^* z_T^* = (W_T^*)^\top x - \lambda \mathbf{1}.$$

846 Under the assumption that there is no active off-cluster contribution into T ,

$$(W_T^*)^\top x = G_{TT}^* f_T.$$

847 Therefore

$$z_T^* = f_T - \lambda (G_{TT}^*)^{-1} \mathbf{1}.$$

848 For an equiangular block,

$$G_{TT}^* = (1 - \rho)I + \rho \mathbf{1}\mathbf{1}^\top.$$

849 The vector $\mathbf{1}$ is an eigenvector with eigenvalue

$$1 + (k-1)\rho.$$

850 Thus

$$(G_{TT}^*)^{-1} \mathbf{1} = \frac{1}{1 + (k-1)\rho} \mathbf{1},$$

851 and hence

$$z_i^* = f_i - \frac{\lambda}{1 + (k-1)\rho} \quad \forall i \in T.$$

852

□

Proposition G.5 (simplex cancellation restores direction recovery). *Assume*

$$\{w_1^*, \dots, w_m^*\}$$

forms a regular simplex, so

$$\sum_{i=1}^m w_i^* = 0.$$

Fix $j \in [m]$, and let $a(x) \geq 0$ be a learned activation satisfying

853
$$\mathbf{1}_{\{a>0\}} = \mathbf{1}_{\{f_j>0\}} \quad \text{a.s.}$$

Assume the off-target mixed moments are exchangeable:

$$\mathbb{E}[f_i a] = \beta \quad \forall i \neq j, \quad \mathbb{E}[f_j a] = \gamma.$$

Then

$$\mu(a) := \mathbb{E}[x a] = (\gamma - \beta) w_j^*.$$

In particular, if $\gamma \neq \beta$, then the decoder moment is parallel to w_j^* ; if $\gamma > \beta$, it lies on the positive ray $\mathbb{R}_+ w_j^*$.

854 *Proof.* By definition,

$$\mu(a) = \mathbb{E}[xa] = \sum_{i=1}^m \mathbb{E}[f_i a] w_i^*.$$

855 Using the exchangeability assumption,

$$\mu(a) = \gamma w_j^* + \beta \sum_{i \neq j} w_i^*.$$

856 Since the simplex relation gives

$$\sum_{i \neq j} w_i^* = -w_j^*,$$

857 we obtain

$$\mu(a) = (\gamma - \beta) w_j^*.$$

858

□

Remark G.6 (symmetry and multiplicity). If a finite orthogonal group $\Gamma \subset O(d)$ preserves both the learned teacher dictionary and the distribution of f , then the induced residual law is Γ -invariant. Therefore

$$859 \quad \mathcal{U}_\lambda(gv; R) = \mathcal{U}_\lambda(v; R) \quad \forall g \in \Gamma.$$

Every local or global maximizer of the oracle utility field then comes with its symmetry orbit. Symmetric learned geometries can restore support-to-direction alignment while simultaneously increasing multiplicity.

860 H Amortized bridge proofs and stress tests

861 The bridge between SAE gradient dynamics, oracle one-slot utility, and one-ReLU realizability has
862 two sides. First, decoder updates follow a gated residual moment and one encoder row tries to realize
863 the oracle gate. Second, amortization gaps and pair-only worlds can change the ranking of directions.

864 The bridge from oracle sparse coding to SAE slots turns on whether the encoder can realize the oracle gate; the stress tests show when amortization or pair-only structure changes the recovery picture.

865 **Deferred decoder-gradient statement.** Proposition below is a characterization of moment-level
866 condition of the decoder gate.

Proposition H.1 (decoder projected-gradient moment). *In the normalized untied SAE model, with activations fixed as functions of the encoder input and with $\|d_r\| = 1$, the sphere-projected decoder gradient flow for slot r is proportional to*

$$867 \quad P_{d_r^\perp} \mathbb{E}[R_{-r} z_r].$$

868 H.1 Proof of Theorem 5.1 and Proposition H.1

869 Write the sample loss as

$$\ell(v; D, E, b_{\text{enc}}) := \frac{1}{2} \|\nu - Dz\|_2^2 + \lambda \mathbf{1}^\top z, \quad z = \sigma(E\nu + b_{\text{enc}}),$$

870 with residual $r = \nu - Dz$. Since $\partial r / \partial d_r = -z_r$ and z_r does not depend on d_r in the untied model,

$$\nabla_{d_r} \ell = -r z_r.$$

871 Taking expectations gives

$$\nabla_{d_r} \mathcal{L}_\lambda = -\mathbb{E}[r z_r].$$

872 Next differentiate through the gate $z_r = (e_r^\top \nu + b_r)_+$. Under the no-boundary-mass assumption
 873 $\mathbb{P}(e_r^\top \nu + b_r = 0) = 0$, we may differentiate almost surely with

$$\frac{\partial z_r}{\partial e_r} = m_r \nu, \quad \frac{\partial z_r}{\partial b_r} = m_r, \quad m_r = \mathbf{1}_{\{e_r^\top \nu + b_r > 0\}}.$$

874 Because only z_r depends on (e_r, b_r) ,

$$\frac{\partial \ell}{\partial z_r} = -d_r^\top r + \lambda.$$

875 Therefore

$$\nabla_{e_r} \ell = (\lambda - d_r^\top r) m_r \nu, \quad \partial_{b_r} \ell = (\lambda - d_r^\top r) m_r.$$

876 Taking expectations gives the displayed formulas.

877 Since

$$R_{-r} = \nu - \sum_{s \neq r} d_s z_s = r + d_r z_r,$$

878 and $\|d_r\| = 1$,

$$d_r^\top r = d_r^\top R_{-r} - z_r.$$

879 Hence if

$$y_r(\nu) := d_r^\top R_{-r}(\nu) - \lambda,$$

880 then

$$\lambda - d_r^\top r = z_r - y_r,$$

881 which yields the equivalent encoder-gradient formulas.

882 For the block decomposition, expand

$$\frac{1}{2} \|R_{-r} - d_r z_r\|_2^2 + \lambda z_r = \frac{1}{2} \|R_{-r}\|_2^2 - z_r d_r^\top R_{-r} + \frac{1}{2} z_r^2 + \lambda z_r.$$

883 Using $y_r = d_r^\top R_{-r} - \lambda$, this becomes

$$\frac{1}{2} \|R_{-r}\|_2^2 - z_r y_r + \frac{1}{2} z_r^2 = \frac{1}{2} \|R_{-r}\|_2^2 - \frac{1}{2} y_r^2 + \frac{1}{2} (z_r - y_r)^2.$$

884 Taking expectations shows that, for fixed R_{-r} , blockwise optimization over (e_r, b_r) is exactly

$$\min_{e_r, b_r} \frac{1}{2} \mathbb{E}[(z_r - y_r)^2].$$

885 Finally, the decoder is constrained to the sphere $\|d_r\| = 1$. The ambient negative gradient direction is
 886 $\mathbb{E}[r z_r]$, so the projected gradient flow is

$$\dot{d}_r = P_{d_r^\perp} \mathbb{E}[r z_r].$$

887 Since $P_{d_r^\perp}(d_r z_r^2) = 0$ pointwise and $R_{-r} = r + d_r z_r$, this is equivalently

$$\dot{d}_r = P_{d_r^\perp} \mathbb{E}[R_{-r} z_r].$$

Theorem H.2 (Oracle utility as the sample-wise activation envelope). *Let $\mathcal{G} := L_+^2(\Omega)$ denote arbitrary nonnegative square-integrable gates. Assume $R \in L^2$, and let $d \in S^{d-1}$. Then $\mathcal{U}_\lambda(d; R) = \sup_{z \in \mathcal{G}} \Gamma_\lambda(d, z; R)$, and the unique maximizer in \mathcal{G} is $g_d^* = (d^\top R - \lambda)_+$. Moreover, for every $z \in \mathcal{G}$,*

888

$$\mathcal{U}_\lambda(d; R) - \Gamma_\lambda(d, z; R) = \mathbb{E} \left[\frac{1}{2} (z - g_d^*)^2 + z(\lambda - d^\top R)_+ \right].$$

Consequently, $\mathcal{A}_\lambda(d; R, \nu) \leq \mathcal{U}_\lambda(d; R)$.

889 **H.2 Proof of Theorem H.2**

890 For fixed d and sample value of (R, ν) , optimize pointwise over nonnegative scalar z :

$$q := d^\top R - \lambda, \quad \sup_{z \geq 0} \left\{ zq - \frac{1}{2}z^2 \right\} = \frac{1}{2}q_+^2, \quad z^* = q_+.$$

891 Thus the pointwise maximizer is

$$g_d^* = (d^\top R - \lambda)_+.$$

892 Moreover, for any $z \geq 0$,

$$\frac{1}{2}q_+^2 - \left(zq - \frac{1}{2}z^2 \right) = \frac{1}{2}(z - q_+)^2 + z(-q)_+.$$

893 Substituting $q = d^\top R - \lambda$ and taking expectations gives the envelope formula and the exact gap
894 identity.

895 **H.3 One-ReLU realizability criterion and closure remark**

896 If $g : \mathbb{R}^d \rightarrow \mathbb{R}_+$ is continuous and piecewise affine, exact one-ReLU realizability means

$$g(\nu) = (e^\top \nu + b)_+.$$

897 Then $\{g > 0\}$ must be a halfspace and g must agree there with one affine map. Conversely, any
898 nonnegative function with those two properties is a single affine-ReLU.

899 Strictly speaking, equality $\mathcal{A}_\lambda(d; R, \nu) = \mathcal{U}_\lambda(d; R)$ means that the oracle gate lies in the L^2 -closure
900 of the one-ReLU class. Under optimizer attainment, this reduces to exact one-ReLU realization as
901 stated in Theorem 5.2.

902 **H.4 Proof of Theorem 5.2**

903 *Proof.* The identity in the theorem implies that any sequence $z_n \in \mathcal{A}_\nu$ with $\Gamma_\lambda(d, z_n; R) \rightarrow \mathcal{U}_\lambda(d; R)$
904 must satisfy $\|z_n - g_d^*\|_{L^2} \rightarrow 0$. Conversely, if $z_n \in \mathcal{A}_\nu$ and $z_n \rightarrow g_d^*$ in L^2 , then the same identity
905 and Cauchy–Schwarz imply $\Gamma_\lambda(d, z_n; R) \rightarrow \mathcal{U}_\lambda(d; R)$, since $g_d^* = 0$ on $\{d^\top R - \lambda < 0\}$ and
906 $(\lambda - d^\top R)_+ \in L^2$. The attained statement follows because equality in the nonnegative gap identity
907 forces $z = g_d^*$ almost surely. \square

Remark H.3 (Different envelope for encoder-measurable unrestricted gates). If instead one
optimizes over all nonnegative $\sigma(\nu)$ -measurable gates, the envelope is generally different.
Writing $Y = d^\top R - \lambda$, one obtains

$$\sup_{z \in L^2_+(\sigma(\nu))} \mathbb{E} \left[zY - \frac{1}{2}z^2 \right] = \frac{1}{2} \mathbb{E} \left[(\mathbb{E}[Y | \nu])_+^2 \right],$$

with optimizer $z = (\mathbb{E}[Y | \nu])_+$. This is not the oracle utility $\mathcal{U}_\lambda(d; R)$ used in this paper. The
gap studied here is therefore the gap between the sample-wise sparse-coding oracle and the
one-ReLU amortized encoder class.

909 **H.5 Subgradient formula for amortized utility**

910 For fixed (e, b) , the map $d \mapsto \Gamma_\lambda(d, e, b; R, \nu)$ is affine with gradient

$$\nabla_d \Gamma_\lambda(d, e, b; R, \nu) = \mathbb{E}[Rz_{e,b}(\nu)].$$

911 Therefore \mathcal{A}_λ is a supremum of affine functionals. Under optimizer attainment and boundedness of
912 the active-gradient set, the standard Danskin–Clarke theorem gives

$$\partial^\circ \mathcal{A}_\lambda(d; R, \nu) = \text{conv} \{ \mathbb{E}[Rz_{e,b}(\nu)] : (e, b) \in \mathcal{M}(d) \}.$$

913 Projecting onto d^\perp yields valid sphere ascent directions.

Proposition H.4 (common-field approximation). *Let*

$$g_r^* = (d_r^\top R_{-r} - \lambda)_+, \quad h_r^* = (d_r^\top r - \lambda)_+.$$

Then

$$\left\| \dot{d}_r - \operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; R_{-r}) \right\| \leq \mathbb{E}[\|R_{-r}\| |z_r - g_r^*|].$$

914 *Moreover,*

$$\left\| \operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; R_{-r}) - \operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; r) \right\| \leq \mathbb{E}[\|r\| |z_r].$$

Consequently,

$$\|\dot{d}_r - \operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; r)\| \leq \mathbb{E}[\|R_{-r}\| |z_r - g_r^*|] + \mathbb{E}[\|r\| |z_r].$$

915 H.6 Proof of Proposition H.4

916 Let $g_r^* = (d_r^\top R_{-r} - \lambda)_+$. By [Theorem H.2](#), the oracle one-slot field for the leave-one-out residual
917 satisfies

$$\operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; R_{-r}) = P_{d_r^\perp} \mathbb{E}[R_{-r} g_r^*].$$

918 By [Proposition H.1](#), the actual decoder flow is

$$\dot{d}_r = P_{d_r^\perp} \mathbb{E}[R_{-r} z_r].$$

919 Subtracting yields

$$\dot{d}_r - \operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; R_{-r}) = P_{d_r^\perp} \mathbb{E}[R_{-r} (z_r - g_r^*)],$$

920 and therefore

$$\left\| \dot{d}_r - \operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; R_{-r}) \right\| \leq \mathbb{E}[\|R_{-r}\| |z_r - g_r^*|],$$

921 which is the amortization term.

922 For the common-field comparison, define

$$h_r^* := (d_r^\top r - \lambda)_+.$$

923 Then

$$\operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; r) = P_{d_r^\perp} \mathbb{E}[r h_r^*].$$

924 Hence

$$\operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; R_{-r}) - \operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; r) = P_{d_r^\perp} \mathbb{E}[R_{-r} g_r^* - r h_r^*].$$

925 Using $R_{-r} = r + d_r z_r$ and $P_{d_r^\perp}(d_r z_r g_r^*) = 0$, this becomes

$$P_{d_r^\perp} \mathbb{E}[r (g_r^* - h_r^*)].$$

926 Since $t \mapsto (t - \lambda)_+$ is 1-Lipschitz and

$$d_r^\top R_{-r} - d_r^\top r = z_r,$$

927 we have

$$|g_r^* - h_r^*| \leq z_r.$$

928 Therefore

$$\left\| \operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; R_{-r}) - \operatorname{grad}_{S^{d-1}} \mathcal{U}_\lambda(d_r; r) \right\| \leq \mathbb{E}[\|r\| |z_r].$$

929 Combining this bound with the amortization term proves the result.

930 **H.7 Amortization-induced ranking reversal and pair-only worlds**

931 This stress-test subsection proves the amortization-rescue claims used in Section 5. The deferred
 932 encoder-measurable envelope remark above clarifies what is and is not meant by oracle utility. The
 933 main points are: an amortized SAE can reverse an oracle ranking only through a direction-dependent
 934 amortization gap; this reversal is impossible when the residual map is affine in the encoder input; and
 935 pair-only worlds provide a simple case where both one-slot and exact multi-slot sparse coding favor
 936 non-ground-truth directions.

Proposition H.5 (Amortization-induced ranking reversal). *For a unit direction d , define*

$$\text{Gap}_\lambda(d; R, \nu) := \mathcal{U}_\lambda(d; R) - \mathcal{A}_\lambda(d; R, \nu) \geq 0.$$

Then for any unit directions u, v ,

937
$$\mathcal{A}_\lambda(u; R, \nu) - \mathcal{A}_\lambda(v; R, \nu) = \mathcal{U}_\lambda(u; R) - \mathcal{U}_\lambda(v; R) + \text{Gap}_\lambda(v; R, \nu) - \text{Gap}_\lambda(u; R, \nu).$$

Consequently, if

$$\mathcal{U}_\lambda(v; R) > \mathcal{U}_\lambda(u; R),$$

then the amortized utility ranks u above v only if

$$\text{Gap}_\lambda(v; R, \nu) - \text{Gap}_\lambda(u; R, \nu) > \mathcal{U}_\lambda(v; R) - \mathcal{U}_\lambda(u; R).$$

938 *Proof.* By definition,

$$\mathcal{A}_\lambda(d; R, \nu) = \mathcal{U}_\lambda(d; R) - \text{Gap}_\lambda(d; R, \nu).$$

939 Applying this identity to u and v , then subtracting, gives

$$\mathcal{A}_\lambda(u) - \mathcal{A}_\lambda(v) = \mathcal{U}_\lambda(u) - \mathcal{U}_\lambda(v) + \text{Gap}_\lambda(v) - \text{Gap}_\lambda(u).$$

940 The strict reversal condition follows immediately. □

Proposition H.6 (No amortized landscape reshaping under affine residual maps). *Assume*

$$R(\nu) = M\nu + r_0 \quad \text{a.s.}$$

941 *for some matrix M and vector r_0 . Then for every unit direction d ,*

$$\mathcal{A}_\lambda(d; R, \nu) = \mathcal{U}_\lambda(d; R).$$

Therefore the amortized and oracle directional landscapes coincide exactly.

942 *Proof.* The oracle gate for direction d is

$$g_d^*(\nu) = (d^\top R(\nu) - \lambda)_+.$$

943 If $R(\nu) = M\nu + r_0$, then

$$g_d^*(\nu) = ((M^\top d)^\top \nu + d^\top r_0 - \lambda)_+.$$

944 This is exactly a single affine-ReLU. By the exact bridge criterion,

$$\mathcal{A}_\lambda(d; R, \nu) = \mathcal{U}_\lambda(d; R).$$

945 Since this holds for every d , the full amortized directional landscape agrees with the oracle landscape. □

946

Remark H.7 (Scope of the affine no-reshaping result). The result applies to the untied one-ReLU encoder class used in the bridge theorem. Tied encoders impose an additional restriction $e_r = d_r$, and therefore may have a nonzero amortization gap even when $R(\nu)$ is affine. Any ranking change in the tied case is due to this extra restriction, not to the generic one-ReLU amortization gap analyzed here.

948 **Pair-only one-slot utility.** Let $u_1, u_2, u_3 \in S^{d-1}$ be orthonormal and let $A > 0$. Suppose the data
 949 distribution is

$$x_{12} = A(u_1 + u_2), \quad x_{13} = A(u_1 + u_3),$$

950 with probability 1/2 each.

951 For the ground-truth direction u_1 ,

$$\langle u_1, x_{12} \rangle = A, \quad \langle u_1, x_{13} \rangle = A.$$

952 Hence

$$\mathcal{U}_\lambda(u_1) = \frac{1}{2} \cdot \frac{1}{2}(A - \lambda)_+^2 + \frac{1}{2} \cdot \frac{1}{2}(A - \lambda)_+^2 = \frac{1}{2}(A - \lambda)_+^2.$$

953 For the pair direction

$$d_{12} := \frac{u_1 + u_2}{\sqrt{2}},$$

954 one has

$$\langle d_{12}, x_{12} \rangle = \sqrt{2}A, \quad \langle d_{12}, x_{13} \rangle = \frac{A}{\sqrt{2}}.$$

955 Therefore

$$\mathcal{U}_\lambda(d_{12}) = \frac{1}{4}(\sqrt{2}A - \lambda)_+^2 + \frac{1}{4}(A/\sqrt{2} - \lambda)_+^2.$$

956 In particular, if

$$A < \lambda < \sqrt{2}A,$$

957 then

$$\mathcal{U}_\lambda(u_1) = 0, \quad \mathcal{U}_\lambda(d_{12}) = \frac{1}{4}(\sqrt{2}A - \lambda)^2 > 0.$$

958 Thus the non-ground-truth pair direction is strictly preferred to the ground-truth direction u_1 . Since
 959 $R(\nu) = \nu = x$ in the one-slot setting, [Proposition H.6](#) implies that the untied one-slot SAE has
 960 exactly the same ranking.

Proposition H.8 (Exact sparse coding prefers pair features in the pair-only model). *In the pair-only model above, define*

$$D_{\text{GT}} = [u_1, u_2, u_3]$$

and

$$d_{12} = \frac{u_1 + u_2}{\sqrt{2}}, \quad d_{13} = \frac{u_1 + u_3}{\sqrt{2}}, \quad D_{\text{pair}} = [d_{12}, d_{13}].$$

961 For an equal-budget comparison, also let $D_{\text{pair}}^+ = [d_{12}, d_{13}, e]$ for any unit vector e . Let

$$\Phi_\lambda(D) = \mathbb{E} \left[\min_{z \geq 0} \frac{1}{2} \|x - Dz\|^2 + \lambda \mathbf{1}^\top z \right].$$

Then

$$\Phi_\lambda(D_{\text{pair}}^+) \leq \Phi_\lambda(D_{\text{pair}}) < \Phi_\lambda(D_{\text{GT}}) \quad \text{for every } 0 < \lambda < \sqrt{2}A,$$

and $\Phi_\lambda(D_{\text{pair}}^+) = \Phi_\lambda(D_{\text{pair}}) = \Phi_\lambda(D_{\text{GT}}) = A^2$ for $\lambda \geq \sqrt{2}A$.

962 *Proof.* First consider D_{GT} . Since u_1, u_2, u_3 are orthonormal, the nonnegative Lasso separates
 963 coordinatewise. For $x_{12} = Au_1 + Au_2$, the optimal code is

$$z_1^* = z_2^* = (A - \lambda)_+, \quad z_3^* = 0.$$

964 The same loss is obtained for x_{13} . Hence

$$\Phi_\lambda(D_{\text{GT}}) = \begin{cases} 2\lambda A - \lambda^2, & 0 \leq \lambda < A, \\ A^2, & \lambda \geq A. \end{cases}$$

965 Indeed, for $0 \leq \lambda < A$, the residual on an active coordinate is λ , giving reconstruction cost λ^2 across
 966 the two active coordinates and sparsity cost $2\lambda(A - \lambda)$, hence total $2\lambda A - \lambda^2$. For $\lambda \geq A$, the zero
 967 code is optimal and the loss is

$$\frac{1}{2}\|x_{12}\|^2 = A^2.$$

968 Now consider D_{pair} . For x_{12} ,

$$x_{12} = \sqrt{2}A d_{12}.$$

969 The candidate code

$$z_{12}^* = (\sqrt{2}A - \lambda)_+, \quad z_{13}^* = 0$$

970 satisfies the KKT conditions. When $0 \leq \lambda < \sqrt{2}A$, the residual is

$$r = x_{12} - d_{12}z_{12}^* = \lambda d_{12}.$$

971 The active KKT condition is $d_{12}^\top r = \lambda$. The inactive KKT inequality is

$$d_{13}^\top r = \lambda d_{13}^\top d_{12} = \frac{\lambda}{2} \leq \lambda.$$

972 Thus the candidate is optimal. The same argument holds for x_{13} by symmetry. Therefore

$$\Phi_\lambda(D_{\text{pair}}) = \begin{cases} \sqrt{2}\lambda A - \frac{1}{2}\lambda^2, & 0 \leq \lambda < \sqrt{2}A, \\ A^2, & \lambda \geq \sqrt{2}A. \end{cases}$$

973 It remains to compare the two expressions. If $0 < \lambda < A$, then

$$\Phi_\lambda(D_{\text{GT}}) - \Phi_\lambda(D_{\text{pair}}) = (2\lambda A - \lambda^2) - \left(\sqrt{2}\lambda A - \frac{1}{2}\lambda^2\right) = \lambda A(2 - \sqrt{2}) - \frac{1}{2}\lambda^2.$$

974 Since $\lambda < A$,

$$\lambda A(2 - \sqrt{2}) - \frac{1}{2}\lambda^2 > \lambda A \left(2 - \sqrt{2} - \frac{1}{2}\right) > 0.$$

975 If $A \leq \lambda < \sqrt{2}A$, then

$$\Phi_\lambda(D_{\text{GT}}) = A^2$$

976 and

$$\Phi_\lambda(D_{\text{pair}}) = \sqrt{2}\lambda A - \frac{1}{2}\lambda^2.$$

977 Thus

$$\Phi_\lambda(D_{\text{GT}}) - \Phi_\lambda(D_{\text{pair}}) = A^2 - \sqrt{2}\lambda A + \frac{1}{2}\lambda^2 = \frac{1}{2}(\lambda - \sqrt{2}A)^2 > 0.$$

978 Since D_{pair}^+ contains the two columns of D_{pair} , its optimal value is no larger than $\Phi_\lambda(D_{\text{pair}})$, giving
 979 the strict equal-budget comparison for $0 < \lambda < \sqrt{2}A$. For $\lambda \geq \sqrt{2}A$, every unit column d satisfies
 980 $d^\top x \leq \|x\| = \sqrt{2}A \leq \lambda$ on both samples, so the zero code is KKT-optimal for all three dictionaries
 981 and the loss is A^2 . This proves the claim. \square

Proposition H.9 (Support-faithful pair-only firing need not recover the atom direction). *In the pair-only model, let $a \geq 0$ be a gate satisfying*

$$a(x_{12}) > 0, \quad a(x_{13}) > 0.$$

982 *Then*

$$\mu(a) := \mathbb{E}[xa]$$

is not parallel to u_1 . In particular, a gate that fires on all events containing latent feature 1 does not recover the decoder direction u_1 .

983 *Proof.* Write

$$a_{12} := a(x_{12}) > 0, \quad a_{13} := a(x_{13}) > 0.$$

984 Then

$$\mu(a) = \frac{1}{2}Aa_{12}(u_1 + u_2) + \frac{1}{2}Aa_{13}(u_1 + u_3).$$

985 Equivalently,

$$\mu(a) = \frac{A}{2}(a_{12} + a_{13})u_1 + \frac{A}{2}a_{12}u_2 + \frac{A}{2}a_{13}u_3.$$

986 Since $a_{12}, a_{13} > 0$ and u_1, u_2, u_3 are orthonormal, the u_2 and u_3 components are nonzero. Hence
987 $\mu(a) \not\parallel u_1$. \square

Proposition H.10 (Generic pair-only support contamination). *Assume the residualized atoms $\{\tilde{w}_i : i \notin B\}$ are linearly independent. Fix $j \notin B$, and suppose*

$$\mathbb{P}(E_{\{j\}}(B)) = 0.$$

Let $a \geq 0$, with $\mathbb{E}[\|R_B\|a] < \infty$, satisfy

988
$$a > 0 \quad \text{a.s. on } \{\alpha_j > 0\}.$$

If there exists $i \notin B, i \neq j$, such that

$$\mathbb{P}(\alpha_i > 0, \alpha_j > 0) > 0,$$

then

$$\mu(a; B) \not\parallel \tilde{w}_j.$$

989 *Proof.* Since $a > 0$ a.s. on $\{\alpha_j > 0\}$, and since

$$\mathbb{P}(\alpha_i > 0, \alpha_j > 0) > 0,$$

990 we have

$$\mathbb{E}[\alpha_i a] > 0.$$

991 The gated moment is

$$\mu(a; B) = \sum_{k \notin B} \mathbb{E}[\alpha_k a] \tilde{w}_k.$$

992 By the moment characterization theorem, under linear independence this vector is parallel to \tilde{w}_j only
993 if all coefficients on $\tilde{w}_i, i \neq j$, vanish. The i -th coefficient is strictly positive, so $\mu(a; B) \not\parallel \tilde{w}_j$. \square

Remark H.11 (Special cancellation is an exception). The negative pair-only statements above are generic for linearly independent positive geometry. They do not exclude special cancellation geometries. For example, in a regular simplex with exchangeable off-target gated moments, the off-target contributions can cancel exactly, as in [Proposition G.5](#). Such cases are structured exceptions rather than consequences of vanilla sparse reconstruction.

995 I Residual-field transfer and bias effects

996 The residual-field transfer results connect the positive latent-ray oracle field

$$\mathcal{U}_\lambda(\cdot; R_B)$$

997 to the SAE leave-one-out field

$$\mathcal{U}_\lambda(\cdot; R_{-r}).$$

998 The transfer is not automatic: it holds under an explicit residual-compatibility condition and a strict
999 margin.

1000 When the latent residual field and the SAE leave-one-out residual differ, biases and residual mismatch spend recovery margin.

1001 **I.1 Shifted support-event decomposition**

Proposition I.1 (Shifted support-event decomposition). *Let*

$$x = \sum_{i=1}^F \alpha_i w_i, \quad \alpha_i \geq 0,$$

and let

$$R_B = \pi_B x = \sum_{i \notin B} \alpha_i \tilde{w}_i, \quad \tilde{w}_i := \pi_B w_i,$$

be the projected residual. For a deterministic shift $b \in \mathbb{R}^d$, define

$$R_{B,b} := \pi_B(x - b) = R_B - \beta_B, \quad \beta_B := \pi_B b.$$

1002 Then, for every unit vector v ,

$$\mathcal{U}_\lambda(v; R_{B,b}) = \sum_{T \subseteq [F] \setminus B} \frac{1}{2} \mathbb{E} \left[\left(\sum_{i \in T} \alpha_i \langle v, \tilde{w}_i \rangle - \langle v, \beta_B \rangle - \lambda \right)_+^2 \mathbf{1}_{E_T(B)} \right].$$

Equivalently, the deterministic shift turns the original scalar threshold λ into the direction-dependent threshold

$$\lambda_v^b := \lambda + \langle v, \beta_B \rangle.$$

In particular, the empty support event $T = \emptyset$ contributes

$$\frac{1}{2} \mathbb{E} [(-\langle v, \beta_B \rangle - \lambda)_+^2 \mathbf{1}_{E_\emptyset(B)}].$$

1003 *Proof.* On the event $E_T(B)$,

$$R_B = \sum_{i \in T} \alpha_i \tilde{w}_i.$$

1004 Therefore

$$R_{B,b} = R_B - \beta_B = \sum_{i \in T} \alpha_i \tilde{w}_i - \beta_B,$$

1005 and

$$\langle v, R_{B,b} \rangle - \lambda = \sum_{i \in T} \alpha_i \langle v, \tilde{w}_i \rangle - \langle v, \beta_B \rangle - \lambda.$$

1006 Substituting into

$$\mathcal{U}_\lambda(v; R_{B,b}) = \frac{1}{2} \mathbb{E} [(\langle v, R_{B,b} \rangle - \lambda)_+^2]$$

1007 and summing over the full partition

$$\{E_T(B) : T \subseteq [F] \setminus B\}$$

1008 gives the formula. □

1009 This differs structurally from Theorem 3.2, where the sum may be restricted to nonempty residual
1010 supports because $R_B = 0$ on $E_\emptyset(B)$. After a shift, $R_{B,b} = -\beta_B$ on the empty support event, so
1011 absence events can contribute positive utility whenever

$$-\langle v, \beta_B \rangle > \lambda.$$

1012 **I.2 Uniform perturbation bounds for utility and gradient**

Lemma I.2 (Uniform utility-field perturbation, expanded). *Let $R, S \in L^2(\mathbb{R}^d)$ and $\lambda \geq 0$. Then*

$$\sup_{\|v\|=1} |\mathcal{U}_\lambda(v; R) - \mathcal{U}_\lambda(v; S)| \leq \frac{1}{2} \|R - S\|_{L^2} (\|R\|_{L^2} + \|S\|_{L^2}) =: \eta(R, S).$$

1014 *Proof of Lemma I.2.* For real numbers a, b ,

$$|a_+^2 - b_+^2| = |a_+ - b_+|(a_+ + b_+) \leq |a - b|(a_+ + b_+).$$

1015 Set

$$a = \langle v, R \rangle - \lambda, \quad b = \langle v, S \rangle - \lambda.$$

1016 Since $\lambda \geq 0$,

$$a_+ \leq |\langle v, R \rangle|, \quad b_+ \leq |\langle v, S \rangle|.$$

1017 Therefore

$$|\mathcal{U}_\lambda(v; R) - \mathcal{U}_\lambda(v; S)| \leq \frac{1}{2} \mathbb{E}[|\langle v, R - S \rangle| (|\langle v, R \rangle| + |\langle v, S \rangle|)].$$

1018 By Cauchy–Schwarz,

$$\mathbb{E}[|\langle v, R - S \rangle| |\langle v, R \rangle|] \leq \|\langle v, R - S \rangle\|_{L^2} \|\langle v, R \rangle\|_{L^2} \leq \|R - S\|_{L^2} \|R\|_{L^2},$$

1019 where the final inequality uses $\|v\| = 1$. Similarly,

$$\mathbb{E}[|\langle v, R - S \rangle| |\langle v, S \rangle|] \leq \|R - S\|_{L^2} \|S\|_{L^2}.$$

1020 Combining the two estimates gives the stated uniform bound. \square

Lemma I.3 (Gradient-field perturbation). *Let $R, S \in L^2(\mathbb{R}^d)$. Then, for every unit vector v ,*

$$\|\nabla_v \mathcal{U}_\lambda(v; R) - \nabla_v \mathcal{U}_\lambda(v; S)\| \leq (\|R\|_{L^2} + \|S\|_{L^2}) \|R - S\|_{L^2}.$$

The same bound holds after projecting onto v^\perp .

1022 *Proof of Lemma I.3.* Using

$$\nabla_v \mathcal{U}_\lambda(v; R) = \mathbb{E}[(\langle v, R \rangle - \lambda)_+ R],$$

1023 write

$$\nabla_v \mathcal{U}_\lambda(v; R) - \nabla_v \mathcal{U}_\lambda(v; S) = \mathbb{E}[(\langle v, R \rangle - \lambda)_+ (R - S)] + \mathbb{E}[(\langle v, R \rangle - \lambda)_+ - (\langle v, S \rangle - \lambda)_+] S.$$

1024 For the first term,

$$(\langle v, R \rangle - \lambda)_+ \leq |\langle v, R \rangle| \leq \|R\|,$$

1025 so Cauchy–Schwarz gives the bound

$$\|\mathbb{E}[(\langle v, R \rangle - \lambda)_+ (R - S)]\| \leq \|R\|_{L^2} \|R - S\|_{L^2}.$$

1026 For the second term, the positive-part map is 1-Lipschitz, hence

$$|(\langle v, R \rangle - \lambda)_+ - (\langle v, S \rangle - \lambda)_+| \leq |\langle v, R - S \rangle| \leq \|R - S\|.$$

1027 Another Cauchy–Schwarz application gives

$$\|\mathbb{E}[(\langle v, R \rangle - \lambda)_+ - (\langle v, S \rangle - \lambda)_+] S\| \leq \|R - S\|_{L^2} \|S\|_{L^2}.$$

1028 Adding the two estimates proves the claim. Orthogonal projection onto v^\perp cannot increase norm. \square

1029 I.3 Margin transfer theorem

1030 For the SAE leave-one-out residual define

$$\varepsilon_{B, -r} := \|R_{-r} - R_B\|_{L^2} = \left\| P_B x - b_{\text{dec}} - \sum_{s \neq r} d_s z_s \right\|_{L^2},$$

1031 and

$$\eta_{B, -r} := \eta(R_B, R_{-r}) = \frac{1}{2} \|R_{-r} - R_B\|_{L^2} (\|R_B\|_{L^2} + \|R_{-r}\|_{L^2}).$$

Theorem I.4 (Residual-field transfer to amortized SAE utility). *Let $\mathcal{C} \subseteq S^{d-1}$ be a competitor set, and suppose a direction u has latent oracle margin*

$$\gamma := \mathcal{U}_\lambda(u; R_B) - \sup_{v \in \mathcal{C}} \mathcal{U}_\lambda(v; R_B) > 0.$$

Let

$$\text{Gap}_\lambda(u; R_{-r}, \nu) = \mathcal{U}_\lambda(u; R_{-r}) - \mathcal{A}_\lambda(u; R_{-r}, \nu)$$

1032

be the amortization gap of u on the SAE residual. Then

$$\mathcal{A}_\lambda(u; R_{-r}, \nu) - \sup_{v \in \mathcal{C}} \mathcal{A}_\lambda(v; R_{-r}, \nu) \geq \gamma - 2\eta_{B,-r} - \text{Gap}_\lambda(u; R_{-r}, \nu).$$

In particular, if

$$\gamma > 2\eta_{B,-r} + \text{Gap}_\lambda(u; R_{-r}, \nu),$$

then the amortized SAE utility ranks u above every direction in \mathcal{C} .

1033 *Proof of Theorems 5.3 and I.4.* By Lemma I.2,

$$\mathcal{U}_\lambda(u; R_{-r}) \geq \mathcal{U}_\lambda(u; R_B) - \eta_{B,-r},$$

1034 and

$$\sup_{v \in \mathcal{C}} \mathcal{U}_\lambda(v; R_{-r}) \leq \sup_{v \in \mathcal{C}} \mathcal{U}_\lambda(v; R_B) + \eta_{B,-r}.$$

1035 Therefore

$$\mathcal{U}_\lambda(u; R_{-r}) - \sup_{v \in \mathcal{C}} \mathcal{U}_\lambda(v; R_{-r}) \geq \gamma - 2\eta_{B,-r}.$$

1036 By definition,

$$\mathcal{A}_\lambda(u; R_{-r}, \nu) = \mathcal{U}_\lambda(u; R_{-r}) - \text{Gap}_\lambda(u; R_{-r}, \nu).$$

1037 Moreover, by Theorem H.2,

$$\mathcal{A}_\lambda(v; R_{-r}, \nu) \leq \mathcal{U}_\lambda(v; R_{-r}) \quad \forall v.$$

1038 Hence

$$\sup_{v \in \mathcal{C}} \mathcal{A}_\lambda(v; R_{-r}, \nu) \leq \sup_{v \in \mathcal{C}} \mathcal{U}_\lambda(v; R_{-r}).$$

1039 Combining the last three displays gives the stated inequality. \square

1040 Applying this theorem with

$$\mathcal{C}_\varepsilon = \{v \in S^{d-1} : \angle(v, u_j) \geq \varepsilon\}$$

1041 gives the transferred singleton-localization statement used in Section 5.1: if

$$\mathcal{U}_\lambda(u_j; R_B) - \sup_{v \in \mathcal{C}_\varepsilon} \mathcal{U}_\lambda(v; R_B) > 2\eta_{B,-r} + \text{Gap}_\lambda(u_j; R_{-r}, \nu),$$

1042 then the amortized SAE utility ranks u_j above every direction at angle at least ε from u_j .

1043 I.4 Transfer of local merge pressure

Corollary I.5 (Transferred local merge pressure). *Let $u \in S^{d-1}$, $h \in u^\perp \cap S^{d-1}$, and suppose the latent oracle gradient satisfies*

$$\rho := \langle \nabla_v \mathcal{U}_\lambda(u; R_B), h \rangle = \mathbb{E}[\langle (u, R_B) - \lambda \rangle_+ \langle h, R_B \rangle] > 0.$$

If

$$\rho > (\|R_B\|_{L^2} + \|R_{-r}\|_{L^2}) \|R_{-r} - R_B\|_{L^2},$$

1044

then

$$\langle \nabla_v \mathcal{U}_\lambda(u; R_{-r}), h \rangle > 0.$$

Consequently, u is not a local maximizer of

$$v \mapsto \mathcal{U}_\lambda(v; R_{-r})$$

on the sphere.

1045 *Proof.* By Lemma I.3,

$$|\langle \nabla_v \mathcal{U}_\lambda(u; R_{-r}) - \nabla_v \mathcal{U}_\lambda(u; R_B), h \rangle| \leq (\|R_B\|_{L^2} + \|R_{-r}\|_{L^2}) \|R_{-r} - R_B\|_{L^2}.$$

1046 The assumed strict inequality implies

$$\langle \nabla_v \mathcal{U}_\lambda(u; R_{-r}), h \rangle > 0.$$

1047 For the great-circle perturbation

$$v_\theta = \cos \theta u + \sin \theta h,$$

1048 we have

$$\left. \frac{d}{d\theta} \mathcal{U}_\lambda(v_\theta; R_{-r}) \right|_{\theta=0} = \langle \nabla_v \mathcal{U}_\lambda(u; R_{-r}), h \rangle > 0.$$

1049 Thus the utility increases along a feasible tangent direction, so u is not a local maximizer on the
1050 sphere. \square

1051 I.5 Special case: bias-only mismatch

1052 Suppose the previously learned slots exactly reconstruct the V_B -component,

$$\sum_{s \neq r} d_s z_s = P_B x.$$

1053 Then

$$P_{V_B^\perp} R_{-r} = R_B - \pi_B b_{\text{dec}}.$$

1054 Consequently,

$$\|P_{V_B^\perp} R_{-r} - R_B\|_{L^2} = \|\pi_B b_{\text{dec}}\|.$$

1055 If $b_{\text{dec}} \in V_B$, then $\pi_B b_{\text{dec}} = 0$, so the projected residual seen by directions in V_B^\perp coincides with
1056 R_B , up to the amortization gap.

1057 J Direction–support separation proofs

1058 Decoder direction is a statement about a gated first moment, whereas support recovery is a statement
1059 about the event on which the gate fires. The proofs below show why the two coincide only under
1060 extra fixed-direction and exact-support assumptions.

1061 Direction recovery and support recovery separate through exact counterexamples; they coincide under rank-one event support.

1062 J.1 Proof of Theorem 6.1

1063 Expand

$$\mu(a; B) = \sum_{i \notin B} \mathbb{E}[\alpha_i a] \tilde{w}_i.$$

1064 Because $\mu(a; B) \neq 0$ and $\tilde{w}_j \neq 0$, direction alignment means $\mu(a; B) \in \text{span}(\tilde{w}_j)$. This holds if
1065 and only if the sum of the off- j contributions lies in $\text{span}(\tilde{w}_j)$, which proves the first claim. Under
1066 linear independence of the residualized atoms, all coefficients except possibly the j -th must vanish.
1067 Since $\alpha_i, a \geq 0$, the zero-moment condition $\mathbb{E}[\alpha_i a] = 0$ is equivalent to $\alpha_i a = 0$ almost surely.

1068 **Proposition J.1** (Direction without support). *Let $x = \alpha_j w_j$ with $\mathbb{E}[\alpha_j] < \infty$. For any $G \subseteq \{\alpha_j > 0\}$ with positive conditional probability, the gate $a = \mathbf{1}_G$ gives $\mu(a) \parallel w_j$. If the conditional law admits positive-probability subsets of arbitrarily small mass, then support recall can be made arbitrarily low while maintaining exact direction recovery.*

1069 **Proposition J.2** (Support without direction). *Let $x = \alpha_1 w_1 + \alpha_2 w_2$ with $w_1 \perp w_2$ and both support events $\{1\}$ and $\{1, 2\}$ having positive probability. Set $a = \mathbf{1}_{\{\alpha_1 > 0\}}$, which recovers the presence event of atom 1 exactly. Then $\mu(a) = \mathbb{E}[\alpha_1 a] w_1 + \mathbb{E}[\alpha_2 a] w_2$; whenever $\mathbb{E}[\alpha_2 a] > 0$ the decoder direction is rotated away from w_1 .*

1070 **J.2 Proofs of Propositions J.1 and J.2**

1071 For Proposition J.1, in the one-atom world one has

$$\mu(a) = \mathbb{E}[\alpha_j a] w_j,$$

1072 so the decoder direction is exactly w_j . The support recall is $\mathbb{P}(G \mid \{\alpha_j > 0\})$; if the conditional law
1073 admits positive-probability subsets of arbitrarily small mass, this recall can be made arbitrarily low
1074 while preserving direction alignment.

1075 For Proposition J.2, the activation recovers the support event $\{\alpha_1 > 0\}$ exactly by construction, but

$$\mu(a) = \mathbb{E}[\alpha_1 a] w_1 + \mathbb{E}[\alpha_2 a] w_2,$$

1076 and the second coefficient is positive whenever the event $\{1, 2\}$ has positive probability and α_2 is
1077 positive there.

1078 **J.3 Proof of Theorem 6.2**

1079 Since $a = 0$ almost surely on A^c and $R_B = \xi u$ almost surely on A ,

$$\mu(a; B) = \mathbb{E}[R_B a] = \mathbb{E}[\xi a \mathbf{1}_A] u.$$

1080 The scalar prefactor is positive by assumption, so the normalized moment direction is exactly u .

1081 **K Sparsity, frequency matching, and support-event geometry**

1082 The support-event view makes sparsity matching and marginal-frequency matching precise but
1083 limited. Expected sparsity and marginal frequencies are low-order summaries of the latent support
1084 law. The sparse-reconstruction objective depends on the full event law $\mathbb{P}(S = T)$, together with the
1085 amplitudes and geometry on each support event.

1086 Matching global sparsity or marginal frequencies is not enough. The missing object is the full joint support-event law, especially pair and burst events that can dominate thresholded utility.

1087 **K.1 Cardinality matching and support recovery**

1088 Let

$$Z_i := \mathbf{1}_{\{\alpha_i > 0\}}, \quad S := \{i : Z_i = 1\}, \quad N := |S| = \sum_{i=1}^F Z_i.$$

1089 Let a learned code have active set

$$\widehat{S}(x) := \{r : z_r(x) > 0\}, \quad \widehat{N}(x) := |\widehat{S}(x)|.$$

1090 The ground-truth expected sparsity is

$$s_* := \mathbb{E}N = \sum_{i=1}^F p_i, \quad p_i := \mathbb{P}(Z_i = 1).$$

1091 Support recovery, samplewise sparsity recovery, and expected sparsity matching satisfy

$$\widehat{S} = S \text{ a.s.} \implies \widehat{N} = N \text{ a.s.} \implies \mathbb{E}\widehat{N} = \mathbb{E}N.$$

1092 The converses fail in general.

Lemma K.1 (Cardinality mismatch lower-bounds support error). *For any two finite sets A, B ,*

$$|A \Delta B| \geq ||A| - |B||.$$

Consequently, if a learned code has fixed samplewise sparsity $\widehat{N} = k$, then

$$\mathbb{E}|\widehat{S} \Delta S| \geq \mathbb{E}|N - k|.$$

In particular, exact support recovery under a deterministic k -sparse code requires

$$N = k \quad \text{a.s.}$$

1094 *Proof.* The identity

$$|A\Delta B| = |A| + |B| - 2|A \cap B|$$

1095 and the inequalities $|A \cap B| \leq |A|$, $|A \cap B| \leq |B|$ imply

$$|A\Delta B| \geq |A| - |B|, \quad |A\Delta B| \geq |B| - |A|.$$

1096 Taking the maximum gives the deterministic inequality. Applying it samplewise gives the expectation
1097 bound. If $\widehat{S} = S$ a.s. and $\widehat{N} = k$ a.s., then $N = k$ a.s. \square

1098 Thus expected sparsity matching fixes only a first moment, and fixed samplewise sparsity fixes only a
1099 cardinality. Neither specifies which support event occurred.

1100 **K.2 Thresholding under the correct orthogonal dictionary**

1101 Even if the dictionary directions are correct, the ℓ_1 threshold can alter support semantics.

Proposition K.2 (Amplitude thresholding under an orthogonal teacher dictionary). *Assume w_1, \dots, w_F are orthonormal and*

$$x = \sum_{i=1}^F \alpha_i w_i, \quad \alpha_i \geq 0.$$

For the teacher dictionary $D = (w_1, \dots, w_F)$, the exact nonnegative ℓ_1 code is

$$z_i^*(x) = (\alpha_i - \lambda)_+.$$

1102 *Hence*

$$\text{supp } z^*(x) = \{i : \alpha_i > \lambda\}.$$

Exact support recovery with the teacher dictionary holds if and only if

$$\mathbb{P}(0 < \alpha_i \leq \lambda) = 0 \quad \text{for every } i.$$

Moreover,

$$\mathbb{E}|\text{supp } z^*(x)| = \sum_{i=1}^F \mathbb{P}(\alpha_i > \lambda) \leq \sum_{i=1}^F \mathbb{P}(\alpha_i > 0) = s_*,$$

with equality if and only if $\mathbb{P}(0 < \alpha_i \leq \lambda) = 0$ for every i .

1103 *Proof.* For an orthonormal dictionary, the objective separates:

$$\frac{1}{2} \left\| x - \sum_i z_i w_i \right\|^2 + \lambda \sum_i z_i = \sum_i \left[\frac{1}{2} (\alpha_i - z_i)^2 + \lambda z_i \right].$$

1104 The nonnegative scalar minimizer is $z_i^* = (\alpha_i - \lambda)_+$. The support and expectation claims follow
1105 immediately. \square

1106 This proposition isolates an amplitude-threshold mismatch. The rest of the appendix concerns a
1107 different issue: even when sparsity statistics are matched, the event geometry can make non-teacher
1108 directions preferable.

1109 **K.3 Marginal frequencies do not determine oracle preference**

1110 The support-event decomposition depends on the full law of S , not only on the marginals $p_i = \mathbb{P}(i \in$
1111 $S)$. The following proposition gives an exact two-feature separation.

Proposition K.3 (Same marginals, different oracle-preferred directions). *Let $u, w \in \mathbb{R}^d$ be orthonormal, let $A > 0$, and set*

$$x = A(Z_1 u + Z_2 w).$$

Assume

$$\mathbb{P}(Z_1 = 1) = \mathbb{P}(Z_2 = 1) = p, \quad \mathbb{P}(Z_1 = Z_2 = 1) = q.$$

Then

$$\mathbb{P}(S = \{1\}) = p - q, \quad \mathbb{P}(S = \{2\}) = p - q, \quad \mathbb{P}(S = \{1, 2\}) = q.$$

Let $m = (u + w)/\sqrt{2}$. For any $\lambda \geq 0$,

$$\mathcal{U}_\lambda(u; x) = \frac{1}{2}p(A - \lambda)_+^2,$$

1112

and

$$\mathcal{U}_\lambda(m; x) = (p - q)\frac{1}{2}(A/\sqrt{2} - \lambda)_+^2 + q\frac{1}{2}(\sqrt{2}A - \lambda)_+^2.$$

In particular, if

$$A/\sqrt{2} < \lambda < A,$$

then

$$\mathcal{U}_\lambda(m; x) = \frac{1}{2}q(\sqrt{2}A - \lambda)^2.$$

For every $0 < p \leq 1/2$, both $q = 0$ and $q = p$ are feasible with the same marginals (p, p) and the same expected sparsity $2p$. For $q = 0$,

$$\mathcal{U}_\lambda(u; x) > \mathcal{U}_\lambda(m; x),$$

whereas for $q = p$,

$$\mathcal{U}_\lambda(m; x) > \mathcal{U}_\lambda(u; x).$$

1113 *Proof.* The formula for $\mathcal{U}_\lambda(u; x)$ follows because $\langle u, x \rangle = A$ exactly when $Z_1 = 1$, which has
 1114 probability p . For m , singleton events have projection $A/\sqrt{2}$, and the joint event has projection $\sqrt{2}A$,
 1115 giving the displayed expression.

1116 If $A/\sqrt{2} < \lambda < A$, singleton contributions to m vanish, while teacher contributions remain. When
 1117 $q = 0$, the merge has zero utility and the teacher has positive utility. When $q = p$,

$$\mathcal{U}_\lambda(m; x) - \mathcal{U}_\lambda(u; x) = \frac{1}{2}p \left[(\sqrt{2}A - \lambda)^2 - (A - \lambda)^2 \right] > 0,$$

1118 because $\sqrt{2}A - \lambda > A - \lambda > 0$. □

1119 Thus marginal frequency matching does not determine recovery. The two laws have the same (p_1, p_2)
 1120 and the same expected sparsity, but different support-event laws.

1121 **K.4 Variance, covariance, and pair-event mass**

1122 Expected sparsity is the first moment of the support size. Pairwise mixed-event mass is controlled by
 1123 second-order structure.

Lemma K.4 (Support-size variance and expected pair mass). Let $N = \sum_{i=1}^F Z_i$, with $p_i = \mathbb{P}(Z_i = 1)$. Then

$$\text{Var}(N) = \sum_{i=1}^F p_i(1 - p_i) + 2 \sum_{i < j} \text{Cov}(Z_i, Z_j).$$

1124

Moreover,

$$\mathbb{E} \binom{N}{2} = \sum_{i < j} \mathbb{P}(Z_i = 1, Z_j = 1) = \frac{1}{2} (\text{Var}(N) + s_*^2 - s_*),$$

where $s_* = \mathbb{E}N$.

1125 *Proof.* The variance identity follows from expanding $\text{Var}(\sum_i Z_i)$. For the pair count,

$$\binom{N}{2} = \frac{1}{2} N(N - 1).$$

1126 Since

$$\mathbb{E}[N(N - 1)] = \mathbb{E}[N^2] - \mathbb{E}[N] = \text{Var}(N) + s_*^2 - s_*,$$

1127 the formula follows. □

1128 For fixed expected sparsity s_* , larger $\text{Var}(N)$ implies larger expected pair co-activation mass. Positive
 1129 covariance produces bursty co-activation and increases mixed-event mass. Negative covariance can
 1130 stabilize the cardinality. However, variance is still only a summary: it does not identify which pairs
 1131 co-activate, their geometry, their amplitudes, or whether their above-threshold contributions dominate
 1132 singleton events. Conversely, zero variance does not imply singleton dominance: if $N = 2$ almost
 1133 surely, then $\text{Var}(N) = 0$ but $\mathbb{E} \binom{N}{2} = 1$.

1134 K.5 Power-law marginal frequencies

1135 Assume independent Bernoulli support indicators with

$$p_i = Ci^{-\alpha}, \quad 0 < C \leq 1, \quad \alpha > 0, \quad i = 1, \dots, F.$$

1136 Let

$$N_F := \sum_{i=1}^F Z_i, \quad s_F := \mathbb{E}N_F = CH_{F,\alpha}, \quad \rho_F := \frac{s_F}{F},$$

1137 where

$$H_{F,\alpha} := \sum_{i=1}^F i^{-\alpha}.$$

Proposition K.5 (Power-law sparsity asymptotics). As $F \rightarrow \infty$,

$$s_F \rightarrow C\zeta(\alpha), \quad \rho_F \sim \frac{C\zeta(\alpha)}{F} \rightarrow 0, \quad \text{if } \alpha > 1,$$

1138

$$s_F \sim C \log F, \quad \rho_F \sim \frac{C \log F}{F} \rightarrow 0, \quad \text{if } \alpha = 1,$$

and

$$s_F \sim \frac{C}{1 - \alpha} F^{1 - \alpha}, \quad \rho_F \sim \frac{C}{1 - \alpha} F^{-\alpha} \rightarrow 0, \quad \text{if } 0 < \alpha < 1.$$

Thus $\rho_F \rightarrow 0$ for every $\alpha > 0$, but $s_F \rightarrow \infty$ whenever $\alpha \leq 1$.

1139 *Proof.* These are the standard integral-test asymptotics for generalized harmonic sums. □

1140 Under independence,

$$\text{Var}(N_F) = \sum_{i=1}^F p_i(1-p_i) = CH_{F,\alpha} - C^2 H_{F,2\alpha}.$$

Proposition K.6 (Power-law variance and singleton disappearance). *Under the independent power-law model above, if $\alpha \leq 1$, then*

$$\frac{\text{Var}(N_F)}{s_F^2} \rightarrow 0 \quad \text{and hence} \quad \frac{N_F}{s_F} \rightarrow 1 \quad \text{in probability.}$$

1141 *In particular,*

$$\mathbb{P}(N_F = 1) \rightarrow 0.$$

Moreover,

$$\mathbb{E} \binom{N_F}{2} = \frac{1}{2} (s_F^2 - C^2 H_{F,2\alpha}) \sim \frac{1}{2} s_F^2 \quad \text{for } \alpha \leq 1.$$

1142 *Proof.* For $0 < \alpha < 1$, $s_F \asymp F^{1-\alpha}$. Also

$$H_{F,2\alpha} = \begin{cases} O(F^{1-2\alpha}), & 0 < \alpha < 1/2, \\ O(\log F), & \alpha = 1/2, \\ O(1), & 1/2 < \alpha < 1. \end{cases}$$

1143 Thus $H_{F,2\alpha} = o(s_F^2)$, and

$$\text{Var}(N_F) = CH_{F,\alpha} - C^2 H_{F,2\alpha} \sim s_F.$$

1144 For $\alpha = 1$, $s_F \sim C \log F$ and $H_{F,2} \rightarrow \zeta(2)$, so again $\text{Var}(N_F)/s_F^2 \rightarrow 0$. Chebyshev's inequality
1145 gives $N_F/s_F \rightarrow 1$ in probability. Since $s_F \rightarrow \infty$, the event $\{N_F = 1\}$ is contained in

$$\left\{ \left| \frac{N_F}{s_F} - 1 \right| \geq 1 - \frac{1}{s_F} \right\},$$

1146 whose probability tends to zero.

1147 Finally, by independence,

$$\mathbb{E} \binom{N_F}{2} = \sum_{i < j} p_i p_j = \frac{1}{2} \left[\left(\sum_i p_i \right)^2 - \sum_i p_i^2 \right] = \frac{1}{2} (s_F^2 - C^2 H_{F,2\alpha}).$$

1148 The same comparison gives the asymptotic equivalence. □

1149 Thus, for $\alpha \leq 1$,

$$\rho_F \rightarrow 0, \quad s_F \rightarrow \infty, \quad \mathbb{P}(N_F = 1) \rightarrow 0, \quad \mathbb{E} \binom{N_F}{2} \sim \frac{1}{2} s_F^2.$$

1150 The representation is sparse as a fraction of all possible features, but it is not singleton-dominated.

1151 **K.6 Load transition at the hidden dimension**

1152 Let d_F be the hidden-state dimension and define the active-load ratio

$$\chi_F := \frac{s_F}{d_F}.$$

1153 This compares the expected number of active features with the ambient dimension.

1154 If $d_F = d$ is fixed, then for $\alpha = 1$,

$$\chi_F \sim \frac{C \log F}{d},$$

1155 so the scale $s_F \approx d$ occurs at

$$F_c \approx \exp(d/C).$$

1156 For $0 < \alpha < 1$,

$$\chi_F \sim \frac{C}{(1-\alpha)d} F^{1-\alpha},$$

1157 so the scale $s_F \approx d$ occurs at

$$F_c \approx \left(\frac{(1-\alpha)d}{C} \right)^{1/(1-\alpha)}.$$

1158 For $\alpha > 1$, $s_F \rightarrow C\zeta(\alpha)$, so there is no F -driven transition; the limiting load is

$$\chi_\infty = \frac{C\zeta(\alpha)}{d}.$$

1159 More generally, if

$$d_F \asymp DF^\beta$$

1160 with $D > 0$, then for $0 < \alpha < 1$,

$$\chi_F \asymp \frac{C}{(1-\alpha)D} F^{1-\alpha-\beta}.$$

1161 Thus the load transition occurs at

$$\beta = 1 - \alpha.$$

1162 If $\beta > 1 - \alpha$, then $\chi_F \rightarrow 0$; if $\beta < 1 - \alpha$, then $\chi_F \rightarrow \infty$.

1163 This load transition is distinct from the sparsity-ratio limit. One can have

$$\rho_F = s_F/F \rightarrow 0$$

1164 while simultaneously

$$s_F/d_F \rightarrow \infty.$$

1165 Thus a representation may be sparse relative to the number of possible features but overloaded relative
1166 to the observed dimension.

1167 **K.7 Pairwise-only supports**

1168 We now analyze a regime in which sparsity matching is exact but singleton dominance is impossible.

1169 Assume

$$|S| = 2 \quad \text{almost surely.}$$

1170 Let

$$\mathbb{P}(S = \{i, j\}) = \pi_{ij}, \quad \sum_{i < j} \pi_{ij} = 1.$$

1171 Then

$$s_\star = 2, \quad \text{Var}(|S|) = 0, \quad s_\star/F = 2/F \rightarrow 0.$$

1172 Nevertheless,

$$\mathbb{P}(S = \{j\}) = 0 \quad \text{for every } j.$$

1173 Therefore the singleton margin used in Theorem 4.1 is zero for every teacher feature.

Proposition K.7 (Pairwise-only supports create local merge pressure). *Assume w_1, \dots, w_F are orthonormal unit vectors. On the pair event $S = \{i, j\}$, suppose*

$$R = \alpha_i w_i + \alpha_j w_j, \quad \alpha_i, \alpha_j > 0.$$

1174 *Fix $i \neq j$. If*

$$\mathbb{E}[(\alpha_i - \lambda)_+ \alpha_j \mathbf{1}_{\{S = \{i, j\}\}}] > 0,$$

then w_i is not a local maximizer of $v \mapsto \mathcal{U}_\lambda(v; R)$ on the sphere.

1175 *Proof.* Let

$$v_\theta = \cos \theta w_i + \sin \theta w_j.$$

1176 The sphere directional derivative is

$$\left. \frac{d}{d\theta} \mathcal{U}_\lambda(v_\theta; R) \right|_{\theta=0} = \mathbb{E}[(\langle w_i, R \rangle - \lambda)_+ \langle w_j, R \rangle].$$

1177 By orthonormality, the event $S = \{i, j\}$ contributes

$$\mathbb{E}[(\alpha_i - \lambda)_+ \alpha_j \mathbf{1}_{\{S=\{i,j\}\}}].$$

1178 Other pair events contribute zero to this tangent derivative: if $i \notin S$, then $\langle w_i, R \rangle = 0$, and if $j \notin S$,
1179 then $\langle w_j, R \rangle = 0$. Hence the derivative is positive, so w_i is not a local maximizer. \square

1180 If amplitudes vary on a pair event, then the event need not be rank-one. On $S = \{i, j\}$,

$$R = \alpha_i w_i + \alpha_j w_j.$$

1181 The normalized residual direction depends on the ratio α_j/α_i . Hence the pair event is rank-one if
1182 and only if this ratio is almost surely constant on the event.

1183 For directions in the pair plane,

$$v_\theta = \cos \theta w_i + \sin \theta w_j,$$

1184 the conditional pair-event utility is

$$U_{ij}(\theta) = \frac{1}{2} \pi_{ij} \mathbb{E}[(\alpha_i \cos \theta + \alpha_j \sin \theta - \lambda)_+^2 \mid S = \{i, j\}].$$

1185 At the teacher endpoint,

$$U'_{ij}(0) = \pi_{ij} \mathbb{E}[(\alpha_i - \lambda)_+ \alpha_j \mid S = \{i, j\}].$$

1186 Thus, whenever the pair event has positive above-threshold mass, the best one-slot compressor of the
1187 pair event is pushed into the pair cone rather than remaining at a teacher ray.

1188 Finally, exact two-sparsity does not by itself identify teacher rays even if exact reconstruction is
1189 possible. Suppose $F = 2$ and

$$x = \alpha_1 w_1 + \alpha_2 w_2, \quad \frac{\alpha_2}{\alpha_1} \in [r_-, r_+] \subset (0, \infty) \quad \text{almost surely.}$$

1190 Then the observed residual directions lie in a truncated cone strictly inside $\text{cone}(w_1, w_2)$. Any two
1191 rays whose positive cone contains

$$\{w_1 + r w_2 : r \in [r_-, r_+]\}$$

1192 can reconstruct every sample using at most two nonnegative coefficients. Thus exact two-sparsity
1193 and zero reconstruction error do not identify the teacher rays unless the observed cone exposes its
1194 boundary rays.

1195 **K.8 Summary**

1196 Sparsity matching fixes a cardinality scale. Frequency matching fixes marginal activation rates.
1197 Neither fixes the support-event law. In the oracle sparse-reconstruction objective, recovery is governed
1198 by thresholded eventwise geometry:

$$\mathbb{P}(S = T), \quad \{\alpha_i : i \in T\}, \quad \{\tilde{w}_i : i \in T\}.$$

1199 Correct sparsity can therefore enable recovery in regimes where it coincides with rank-one or
1200 singleton-event dominance. Outside those regimes, the same matched sparsity can coexist with merge
1201 pressure, pair-cone compression, splitting, or direction/support mismatch.

1202 **L What can be certified by an observer?**

1203 If direction recovery does not imply support recovery, then decoder alignment alone cannot certify
 1204 a trained SAE feature, especially in real models with no ground-truth supports. Any population
 1205 certificate that uses only observed activations must be a functional of the observed activation law. The
 1206 observable substitute is residual-direction geometry relative to a known observer-chosen subspace.

1207 The observer diagnostics separate direction drift, poor cap precision, poor recall, intrinsic cap ambiguity,
 gate mismatch, and residual-path sensitivity.

1208 **L.1 Observable atoms, caps, and directional margins**

1209 Let $H \subseteq \mathbb{R}^d$ be a known observer-chosen subspace, for example the span of selected learned decoder
 1210 directions or any other subspace specified without using latent supports. Since H is known, P_{H^\perp}
 1211 is computable, and $R_H := P_{H^\perp} x$ is observable from x . On $\{R_H \neq 0\}$, define $U_H := R_H / \|R_H\|$.
 1212 The exact observable analogue of a rank-one event is the residual-direction atom $A_u(H) := \{R_H \neq$
 1213 $0, U_H = u\}$.

1214 The observer cannot ask whether an SAE feature recovered the hidden latent event, because that event
 1215 is not observed. The replacement question is geometric: is there an observed residual-direction atom
 1216 or cap, and does sparse coding prefer directions near it? This is weaker than latent support recovery,
 1217 but it is an audit target that depends only on observed activations.

1218 **Proposition L.1** (maximal observable rank-one event). *Fix H and $u \in S^{d-1}$. If $A \in \sigma(R_H)$,
 $A \subseteq \{R_H \neq 0\}$, and $R_H = \xi u$ a.s. on A for some $\xi > 0$ on A , then $A \subseteq A_u(H)$ up to null
 sets. Conversely, $A_u(H)$ is itself a rank-one residual event whenever it has positive probability.*

1219 *Proof.* If $A \subseteq \{R_H \neq 0\}$ and $R_H = \xi u$ with $\xi > 0$ on A , then $U_H = u$ almost surely on A , so
 1220 $A \subseteq A_u(H)$ up to null sets. Conversely, if $\mathbb{P}(A_u(H)) > 0$, then on $A_u(H)$ we have

$$R_H = \|R_H\| u \quad \text{a.s. on } A_u(H),$$

1221 so $A_u(H)$ is a rank-one event for $\mathbb{R}_+ u$. □

1222 Thus $A_u(H)$ is the maximal observable event on which the residual lies on the ray $\mathbb{R}_+ u$. Since exact
 1223 atoms may have zero mass in the population, use positive-width caps

$$A_{u,\varepsilon}(H) := \{R_H \neq 0, \angle(U_H, u) \leq \varepsilon\}.$$

1224 These caps are observable relaxations of rank-one residual events.

1225 For $0 < \varepsilon < \phi \leq \pi$, define

$$M_{\text{dir}}(u, \varepsilon, \phi, \lambda; H) := \underline{W}_\lambda(u, \varepsilon; H) - \overline{W}_\lambda(u, \varepsilon, \phi; H) - L(u, \varepsilon; H),$$

1226 where

$$\underline{W}_\lambda(u, \varepsilon; H) := \frac{1}{2} \mathbb{E}[(\|R_H\| \cos \varepsilon - \lambda)_+^2 \mathbf{1}_{A_{u,\varepsilon}(H)}],$$

1227

$$\overline{W}_\lambda(u, \varepsilon, \phi; H) := \frac{1}{2} \mathbb{E}[(\|R_H\| \cos(\phi - \varepsilon) - \lambda)_+^2 \mathbf{1}_{A_{u,\varepsilon}(H)}],$$

1228 and

$$L(u, \varepsilon; H) := \frac{1}{2} \mathbb{E}[\|R_H\|^2 \mathbf{1}_{A_{u,\varepsilon}(H)^c}].$$

1229 The margin M_{dir} is an observer-side version of singleton dominance. The first term lower-bounds
 1230 how much the cap supports u . The second asks how much a direction at least ϕ away could still get
 1231 from the same cap. The third gives every sample outside the cap to the competitor as a conservative
 1232 budget.

1233 **Theorem L.2** (observer-only directional certificate). *If $M_{\text{dir}}(u, \varepsilon, \phi, \lambda; H) > 0$, then every
 global maximizer of $v \mapsto \mathcal{U}_\lambda(v; R_H)$ lies in the spherical ϕ -ball around u .*

1234 *Proof.* On the cap $A_{u,\varepsilon}(H)$, one has

$$\langle u, R_H \rangle \geq \|R_H\| \cos \varepsilon,$$

1235 which gives the lower bound $\underline{W}_\lambda(u, \varepsilon; H)$ for the utility of u . If $\angle(v, u) \geq \phi$, then on the same cap

$$\angle(v, U_H) \geq \phi - \varepsilon,$$

1236 hence

$$\langle v, R_H \rangle \leq \|R_H\| \cos(\phi - \varepsilon),$$

1237 which yields the upper bound $\overline{W}_\lambda(u, \varepsilon, \phi; H)$. On the complement, competitor utility is bounded
1238 above by $L(u, \varepsilon; H)$, while the utility of u remains nonnegative. Subtracting gives the result. \square

1239 L.2 Feature-level audit quantities

1240 An observer can certify an objective-level statement without latent supports: the observed residual
1241 distribution's sparse-coding utility field prefers a direction near u . A trained feature r , with decoder d_r
1242 and activation a_r , can then be audited against the certified cap. Direction alignment is $D_{r,u} := \langle d_r, u \rangle$.
1243 Cap precision and recall are

$$\pi_{r,u} := \mathbb{P}(A_{u,\varepsilon}(H) \mid a_r > 0), \quad \rho_{r,u} := \mathbb{P}(a_r > 0 \mid A_{u,\varepsilon}(H)).$$

1244 Intrinsic cap ambiguity is

$$\kappa_{\text{in}}(u, \varepsilon; H) := \frac{\mathbb{E}[\|P_{u^\perp} R_H\|^2 \mathbf{1}_{A_{u,\varepsilon}(H)}]}{\mathbb{E}[\|R_H\|^2 \mathbf{1}_{A_{u,\varepsilon}(H)}]}.$$

1245 The one-slot gate gap is $G_r := \mathbb{E}[(a_r - g_r^*)^2] / (\mathbb{E}[(g_r^*)^2] + \delta_0)$, where $g_r^*(\nu) := (d_r^\top R_{-r}(\nu) - \lambda)_+$
1246 and $\delta_0 > 0$ is a fixed stabilizer. Ratios are interpreted only when their denominators are positive.

1247 These quantities separate observable failure modes. Low $D_{r,u}$ indicates direction drift. Low precision
1248 indicates that the feature fires off the certified cap, as in merging or contamination. Low recall
1249 indicates that the feature covers only part of the certified cap, as in splitting, partial coverage, or
1250 absorption-like false negatives. Large κ_{in} indicates that the cap itself is not close to rank-one, so even
1251 a perfect cap detector would have ambiguous direction semantics. Large G_r indicates that the learned
1252 activation fails to realize the sample-wise oracle gate for its decoder and leave-one-out residual.

Proposition L.3 (observable contamination floor). *Let $A = A_{u,\varepsilon}(H)$, assume $\mathbb{E}[\|R_H\|^2 \mathbf{1}_A] > 0$, and suppose $A = G \sqcup C$, where the clean part satisfies*

$$\|P_{u^\perp} R_H\|^2 \leq \sin^2 \varepsilon_0 \|R_H\|^2 \quad \text{on } G.$$

Define the energy fraction of the contaminating part by

$$\omega_C(u, \varepsilon; H) := \frac{\mathbb{E}[\|R_H\|^2 \mathbf{1}_C]}{\mathbb{E}[\|R_H\|^2 \mathbf{1}_A]}.$$

Then

$$\omega_C(u, \varepsilon; H) \geq \frac{\kappa_{\text{in}}(u, \varepsilon; H) - \sin^2 \varepsilon_0}{1 - \sin^2 \varepsilon_0}.$$

In particular, for an exact rank-one refinement $\varepsilon_0 = 0$,

$$\omega_C(u, \varepsilon; H) \geq \kappa_{\text{in}}(u, \varepsilon; H).$$

1254 *Proof.* On the clean part G ,

$$\|P_{u^\perp} R_H\|^2 \leq \sin^2 \varepsilon_0 \|R_H\|^2,$$

1255 while on C ,

$$\|P_{u^\perp} R_H\|^2 \leq \|R_H\|^2.$$

1256 Taking expectations over $A_{u,\varepsilon}(H)$ gives

$$\mathbb{E}[\|P_{u^\perp} R_H\|^2 \mathbf{1}_A] \leq \sin^2 \varepsilon_0 \mathbb{E}[\|R_H\|^2 \mathbf{1}_G] + \mathbb{E}[\|R_H\|^2 \mathbf{1}_C].$$

1257 Divide by $\mathbb{E}[\|R_H\|^2 \mathbf{1}_A]$ and rearrange. \square

1258 **L.3 Recursive and feature-level refinements**

1259 The observer-only theory suggests a recursive recovery procedure. Start with the residualization
 1260 subspace $H_0 = \{0\}$. At stage t , form the projected residual R_{H_t} , estimate the residual-direction law
 1261 U_{H_t} , identify candidate caps $A_{u,\varepsilon}(H_t)$, certify them using M_{dir} , and add certified directions to H_{t+1} .
 1262 This is the residualized analogue of separable NMF: one searches for anchor directions in projected
 1263 residual space rather than anchor points in data space.

Proposition L.4 (recursive residual-direction recovery). *Suppose there exist directions u_1, \dots, u_m and scales $(\varepsilon_t, \phi_t, \lambda_t)$ such that for each stage t ,*

$$M_{\text{dir}}(u_t, \varepsilon_t, \phi_t, \lambda_t; H_{t-1}) > 0, \quad H_{t-1} := \text{span}\{u_1, \dots, u_{t-1}\}.$$

1264 *Then at every stage t , every global maximizer of the oracle utility field*

$$v \mapsto \mathcal{U}_{\lambda_t}(v; R_{H_{t-1}})$$

lies in the ϕ_t -ball around u_t .

1265 *Proof.* Apply **Theorem L.2** at each stage t with residualization subspace H_{t-1} . The positivity of
 1266 the directional margin implies that every global maximizer of the stage- t oracle field lies inside the
 1267 corresponding ϕ_t -ball around u_t . □

Remark L.5 (global sparsity summaries are not support certificates). There exist two disjoint
 witness families A and C with the same probability mass such that the gates

$$a = \mathbf{1}_A, \quad b = \mathbf{1}_C$$

1268 satisfy

$$\mathbb{E}[a] = \mathbb{E}[b], \quad \text{Var}(a) = \text{Var}(b),$$

while recovering disjoint support families.

1269 *Proof.* Take any two disjoint measurable sets A and C with $\mathbb{P}(A) = \mathbb{P}(C) > 0$. The Bernoulli gates
 1270 $a = \mathbf{1}_A$ and $b = \mathbf{1}_C$ then have identical means and variances, but are supported on disjoint event
 1271 families. □

1272 For the feature-level angle bound, fix a target direction u , a target event A such as an observable
 1273 residual-direction cap, and a trained activation $a_r \geq 0$. Define the in-event longitudinal signal and
 1274 the in/out transverse loads by

$$\beta_{r,u}(B) := \mathbb{E}[\langle u, R_B \rangle a_r \mathbf{1}_A],$$

1275

$$\eta_{r,u}^{\text{in}}(B) := \mathbb{E}[\|P_{u^\perp} R_B\| a_r \mathbf{1}_A], \quad \eta_{r,u}^{\text{out}}(B) := \mathbb{E}[\|R_B\| a_r \mathbf{1}_{A^c}].$$

1276 The bound below is meaningful when $\beta_{r,u}(B) > 0$ and the displayed denominator is positive.

Proposition L.6 (feature-level angle bound). *Let*

$$\mu_r(B) := \mathbb{E}[R_B a_r].$$

Then

$$\|P_{u^\perp} \mu_r(B)\| \leq \eta_{r,u}^{\text{in}}(B) + \eta_{r,u}^{\text{out}}(B),$$

and

$$\langle u, \mu_r(B) \rangle \geq \beta_{r,u}(B) - \eta_{r,u}^{\text{out}}(B).$$

Hence, if

1277

$$\beta_{r,u}(B) > \eta_{r,u}^{\text{in}}(B) + 2\eta_{r,u}^{\text{out}}(B),$$

then

$$\tan \angle(\mu_r(B), u) \leq \frac{\eta_{r,u}^{\text{in}}(B) + \eta_{r,u}^{\text{out}}(B)}{\beta_{r,u}(B) - \eta_{r,u}^{\text{out}}(B)}.$$

When the decoder is fit by least squares to the gate, the same bound controls the angle of d_r . At a stationary point of the normalized untied SAE, d_r is aligned with the gated residual moment from [Proposition H.1](#), so the same signal-versus-contamination interpretation continues to apply.

1278 *Proof.* Decompose

$$\mu_r(B) = \mathbb{E}[R_B a_r \mathbf{1}_A] + \mathbb{E}[R_B a_r \mathbf{1}_{A^c}].$$

1279 The first term has longitudinal component at least $\beta_{r,u}(B)$ and transverse norm at most $\eta_{r,u}^{\text{in}}(B)$. The
1280 second term has norm at most $\eta_{r,u}^{\text{out}}(B)$. Therefore

$$\|P_{u^\perp} \mu_r(B)\| \leq \eta_{r,u}^{\text{in}}(B) + \eta_{r,u}^{\text{out}}(B)$$

1281 and

$$\langle u, \mu_r(B) \rangle \geq \beta_{r,u}(B) - \eta_{r,u}^{\text{out}}(B).$$

1282 Taking the ratio yields the tangent bound. □

1283 **Observable contamination-floor lower bound.** For convenience one may define

$$\underline{\omega}(u, \varepsilon; \varepsilon_0, H) := \max \left\{ 0, \frac{\kappa_{\text{in}}(u, \varepsilon; H) - \sin^2 \varepsilon_0}{1 - \sin^2 \varepsilon_0} \right\}.$$

1284 This is the smallest in-cap contamination fraction consistent with an ε_0 -clean refinement. The quantity
1285 is interpreted only when $\mathbb{E}[\|R_H\|^2 \mathbf{1}_{A_{u,\varepsilon}(H)}] > 0$.

1286 **Reporting protocol.** For each certified residual-direction cap, report

$$q := \mathbb{P}(A_{u,\varepsilon}(H)), \quad \kappa_{\text{in}}, \quad M_{\text{dir}}, \quad \underline{\omega}.$$

1287 For each matched feature, report

$$D_{r,u}, \quad \pi_{r,u}, \quad \rho_{r,u}, \quad \eta_{r,u}^{\text{in}}, \quad \eta_{r,u}^{\text{out}}, \quad G_r.$$

1288 This cap-wise audit gives population evidence for faithful recovery, direction-only recovery, intrinsic
1289 cap ambiguity, splitting, merging, absorption, and amortization failure.

1290 **L.4 Failure-mode taxonomy from observer-only quantities**

1291 For completeness, we restate the observer-only taxonomy implied by the quantities above.

1292 **Faithful cap recovery.** A feature matched to a certified cap has high direction score, high cap
1293 precision and recall, low cap contamination, and small gate gap.

1294 **Direction-only recovery.** A feature has high decoder-cap alignment but poor cap precision or recall.
1295 The decoder aligns with the cap direction, but the gate does not fire exactly on the observable cap.

1296 **Intrinsic cap ambiguity.** A cap has positive directional margin but large contamination floor. Poor
1297 support scores in this regime should not be overinterpreted as an SAE failure.

1298 **Merged features.** A feature has no dominant matched cap or draws comparable gate-weighted
1299 signal from several certified caps. The feature is supported by mixed events rather than by a single
1300 rank-one event.

1301 **Splitting.** Several features align with the same residual-direction cap, each with moderate recall,
1302 while their union covers the cap.

1303 **Absorption.** A certified residual-direction cap has positive directional margin, but its residual utility
1304 drops sharply after subtracting learned slots. One natural observable score is

$$\text{Absorb}(u; H) := \mathcal{U}_\lambda(u; R_H) - \mathcal{U}_\lambda\left(u; R_H - \sum_r d_r a_r\right).$$

1305 Large positive absorption indicates that learned features suppress observable residual utility along u ,
1306 consistent with absorption of the cap-supporting event.

1307 **Amortization failure.** A cap is directionally certified, yet no matched feature has small gate gap.

1308 **Residual-path sensitivity.** Small margins or comparable cap utilities flag possible sensitivity to
1309 perturbations; they are not by themselves an SGD seed-instability certificate.

1310 M Synthetic two-stage SAE illustrations

1311 This appendix gives small synthetic illustrations for the realizability and direction/support separations
1312 in the main text. The experiments are not used as evidence for an SGD theory, nor do they validate
1313 the population oracle path. They are sanity checks for three narrower points: a directly observed
1314 learned code is easier to recover than the original GT support law; high decoder alignment need
1315 not imply support recovery; and learned SAE supports can have a different joint support law from
1316 independent GT supports.

1317 **Data and architecture.** We generate data from a positive latent-ray model with $d = 3$ observed
1318 dimensions and $F = 7$ ground-truth features. Supports are independent Bernoulli variables with
1319 probability $p = 0.2$, active magnitudes are independent $\text{Unif}(0, 1)$, and teacher directions are random
1320 unit vectors in \mathbb{R}^3 . Thus the expected GT sparsity is $Fp = 1.4$. Evaluation uses 50,000 fresh
1321 examples.

1322 Each SAE has an untied encoder and decoder. In the row-wise implementation,

$$z = \text{ReLU}((x - b_{\text{dec}})W_{\text{enc}} + b_{\text{enc}}), \quad \hat{x} = zW_{\text{dec}} + b_{\text{dec}},$$

1323 with normalized decoder rows. The objective is

$$\frac{1}{2} \mathbb{E} \|x - \hat{x}\|_2^2 + \lambda \mathbb{E} \|z\|_1.$$

1324 At each training step, decoder gradients are projected to remove radial components, and decoder rows
1325 are renormalized to unit norm after the optimizer update.

1326 **Training protocol.** All SAEs are trained with Adam using learning rate 3×10^{-3} , betas (0.9, 0.999),
1327 $\epsilon = 10^{-8}$, no weight decay, and batch size 2048. SAE1 is trained for up to 4000 steps and SAE2 for
1328 up to 3000 steps. Validation uses a fixed batch of 8192 samples, evaluated every 500 steps. Early
1329 stopping is enabled after 2500 steps; we track EMA validation explained variance and EMA L_1 with
1330 decay 0.95, and stop if both plateau for 20 validation checks, with tolerances 10^{-4} for EV and 10^{-3}
1331 relative improvement for L_1 . This early-stopping criterion did not materially affect the reported runs,
1332 which typically reached the maximum step budget.

1333 For each trained SAE, λ is selected by a small binary search over $[10^{-5}, 10^{-1}]$, with up to two
1334 upper-bound expansions and four geometric midpoint rounds. Each candidate trains a fresh SAE. The

Table 2: Aggregate recovery in the two-stage SAE experiments. Code-student variants train SAE2 on SAE1 activations; output-student variants train SAE2 on SAE1 reconstructed hidden states. Direction-matched support F1 is the main metric for testing whether direction recovery carries the intended activation pattern.

Experiment	Pair	Direction	Direction-matched F1	Support-matched F1	Support-matched direction
Baseline code-student	GT → SAE1	0.953	0.583	0.605	0.952
Baseline code-student	GT → SAE2	0.971	0.600	0.600	0.971
Baseline code-student	SAE1 → SAE2	0.993	0.937	0.937	0.993
Output-student	GT → SAE1	0.953	0.583	0.605	0.952
Output-student	GT → SAE2	0.945	0.514	0.538	0.911
Output-student	SAE1 → SAE2	0.956	0.804	0.804	0.956
Fixed-decoder code-student	GT → SAE1	1.000	0.571	0.577	0.994
Fixed-decoder code-student	GT → SAE2	1.000	0.569	0.575	0.994
Fixed-decoder code-student	SAE1 → SAE2	1.000	0.995	0.995	1.000
Fixed-decoder output-student	GT → SAE1	1.000	0.571	0.577	0.994
Fixed-decoder output-student	GT → SAE2	1.000	0.391	0.518	0.940
Fixed-decoder output-student	SAE1 → SAE2	1.000	0.592	0.711	0.935
Random teacher	GT → SAE1	0.933	0.490	0.557	0.618
Random teacher	GT → SAE2	0.899	0.421	0.518	0.663
Random teacher	SAE1 → SAE2	0.948	0.783	0.818	0.808

1335 selected candidate is the one whose final validation L_0 is closest to the target, breaking ties by higher
 1336 explained variance. SAE1 targets the GT expected sparsity Fp , while SAE2 targets the measured
 1337 teacher L_0 . Default seeds are fixed: toy distribution seed 0, GT dictionary seed 1, SAE1 seed 2,
 1338 SAE2 seed 20, and evaluation seed 4.

1339 **Variants.** We compare five variants. In the *baseline code-student* variant, SAE1 is trained on GT
 1340 hidden states and SAE2 is trained directly on SAE1 activations $z_1(x) = \text{SAE1.encode}(x)$. The
 1341 second-stage target is therefore a nonnegative coordinate code. In the *output-student* variant, SAE2 is
 1342 trained on SAE1 reconstructed hidden states

$$\hat{x}_1(x) = \text{SAE1.decode}(\text{SAE1.encode}(x)),$$

1343 so the teacher code is not given directly. In the *fixed-decoder code-student* variant, SAE1’s decoder is
 1344 fixed to the GT directions, SAE2 is trained on SAE1 activations, and SAE2’s decoder is fixed to the
 1345 identity in SAE1-code space; after composition through SAE1, SAE2’s hidden-space directions are
 1346 GT-aligned. In the *fixed-decoder output-student* variant, SAE1’s decoder is fixed to the GT directions,
 1347 SAE2 is trained on SAE1 reconstructed hidden states, and SAE2’s decoder is fixed to the same GT
 1348 directions in hidden space. In the *random-teacher* variant, SAE1 is left at random initialization and
 1349 SAE2 is trained on the resulting random teacher activations.

1350 **Metrics.** For a source representation A and target representation B , $A \rightarrow B$ means that each source
 1351 feature of A is matched to its best counterpart in B . Direction is the mean best raw cosine between
 1352 source and target decoder directions. Direction-matched F1 first matches by closest decoder direction
 1353 and then computes support F1. Support-matched F1 instead matches by activation support. Unless
 1354 stated otherwise, support F1 in the text means direction-matched F1, since this directly tests whether
 1355 direction recovery carries the intended activation pattern.

1356 All support comparisons use the same underlying evaluation examples. GT supports come from
 1357 the latent generator. SAE1 supports are computed from $z_1(x)$. Code-student SAE2 supports are
 1358 computed from SAE2 applied to $z_1(x)$. Output-student SAE2 supports are computed from SAE2
 1359 applied to $\hat{x}_1(x)$. For SAE1→SAE2 in code-student experiments, directions are compared in SAE1-
 1360 code space. For GT→SAE2 in code-student experiments, SAE2 decoder directions are composed
 1361 back into hidden space through SAE1’s decoder: if D_1 is the SAE1 decoder in hidden space and D_2
 1362 is the SAE2 decoder in SAE1-code space, the composed hidden-space directions are $D_1 D_2$.

1363 The code-student baseline is a positive control. Since SAE2 is trained on SAE1’s nonnegative code,
 1364 the target support representation is directly coordinate-ReLU representable. SAE2 recovers SAE1’s
 1365 learned code much more faithfully than either SAE recovers the original GT supports:

$$F1(\text{SAE1} \rightarrow \text{SAE2}) = 0.937, \quad F1(\text{GT} \rightarrow \text{SAE1}) = 0.583, \quad F1(\text{GT} \rightarrow \text{SAE2}) = 0.600.$$

Table 3: Training and support-correlation diagnostics. “Selection L_0 ” is the final validation L_0 used by the L_1 -selection procedure; “Eval L_0 ” is measured on the 50,000-sample evaluation set used for support correlations. GT supports are nearly independent on the evaluation set, while learned SAE supports have much larger absolute off-diagonal correlations. Approximate L_0 targeting does not determine the support-event law.

Experiment	Representation	L_1	Target L_0	Selection L_0	Eval L_0	EV	Mean abs. off-diag. corr.	Max abs. off-diag. corr.
Baseline code-student	GT	–	–	–	1.395	–	0.003	0.017
Baseline code-student	SAE1	0.237	1.400	1.266	1.260	0.843	0.205	0.561
Baseline code-student	SAE2	0.018	1.273	1.117	1.115	0.995	0.184	0.454
Output-student	GT	–	–	–	1.395	–	0.003	0.017
Output-student	SAE1	0.237	1.400	1.266	1.260	0.843	0.205	0.561
Output-student	SAE2	0.100	1.273	1.039	1.025	0.960	0.161	0.581
Fixed-decoder code-student	GT	–	–	–	1.395	–	0.003	0.017
Fixed-decoder code-student	SAE1	0.487	1.400	1.170	1.158	0.624	0.212	0.561
Fixed-decoder code-student	SAE2	10^{-5}	1.175	1.179	1.166	1.000	0.210	0.561
Fixed-decoder output-student	GT	–	–	–	1.395	–	0.003	0.017
Fixed-decoder output-student	SAE1	0.487	1.400	1.170	1.158	0.624	0.212	0.561
Fixed-decoder output-student	SAE2	0.056	1.175	1.153	1.144	0.968	0.147	0.756
Random teacher	GT	–	–	–	1.395	–	0.003	0.017
Random teacher	SAE1	–	–	3.473	3.468	–	0.386	0.816
Random teacher	SAE2	0.003	3.473	3.303	3.303	0.999	0.249	0.646

1366 The fixed-decoder code-student variant reaches nearly exact recovery, with

$$F1(\text{SAE1} \rightarrow \text{SAE2}) = 0.995.$$

1367 This confirms that the metric and training setup can detect near-perfect recovery when the target code
1368 is directly presented.

1369 The output-student variants are harder because SAE2 is trained on $\hat{x}_1(x)$, not on $z_1(x)$. With learned
1370 decoders, recovery of the SAE1-induced representation is weaker than in the direct code-student case
1371 but remains higher than GT recovery:

$$F1(\text{SAE1} \rightarrow \text{SAE2}) = 0.804, \quad F1(\text{GT} \rightarrow \text{SAE2}) = 0.514.$$

1372 The fixed-decoder output-student ablation is more stringent. Here both SAE1 and SAE2 have the
1373 correct GT decoder directions fixed in hidden space. Nevertheless,

$$F1(\text{GT} \rightarrow \text{SAE2}) = 0.391, \quad F1(\text{SAE1} \rightarrow \text{SAE2}) = 0.592.$$

1374 Thus even when decoder directions are correct by construction, the student does not recover the
1375 teacher or GT support events. The high support-matched directions show that the best support
1376 matches are still geometrically nearby, but the corresponding support overlaps remain far from exact.
1377 Direct code recovery and hidden-state recovery are different tasks.

1378 The fixed-decoder ablations give the cleanest direction/support separation. Decoder directions are
1379 correct by construction, but GT support recovery remains poor. In the code-student fixed-decoder
1380 variant,

$$\text{direction}(\text{GT} \rightarrow \text{SAE2}) = 1.000, \quad F1(\text{GT} \rightarrow \text{SAE2}) = 0.569.$$

1381 In the output-student fixed-decoder variant, the separation is even sharper:

$$\text{direction}(\text{GT} \rightarrow \text{SAE2}) = 1.000, \quad F1(\text{GT} \rightarrow \text{SAE2}) = 0.391.$$

1382 Thus correct decoder directions do not certify recovery of the intended activation pattern: the learned
1383 gates remain geometrically related to the correct rays, but their support events differ substantially.

1384 The fixed-decoder output-student ablation also separates approximate sparsity targeting from support
1385 recovery. SAE2’s final L_0 is 1.153, close to its target 1.175, and its explained variance is 0.968. Yet
1386 its GT support F1 is only 0.391. Thus a model can have correct decoder directions, approximately
1387 matched sparsity, and high reconstruction quality while still using the wrong support-event law.

1388 The correlation check shows the same point at the level of joint support statistics. GT supports
1389 are independent to finite-sample precision, with mean absolute off-diagonal support correlation
1390 0.0034 and maximum 0.0168. Learned supports are substantially correlated. In the baseline, SAE1
1391 has mean 0.2048 and maximum 0.5607, while SAE2 has mean 0.1838 and maximum 0.4544. In

1392 the fixed-decoder output-student variant, SAE2 has mean 0.1473 and maximum 0.7560. Thus the
1393 mismatch is not merely a direction-matching artifact; learned SAE codes can have different joint
1394 support statistics from the GT latent code even when decoder directions are fixed.

1395 These experiments are intentionally limited. They use one synthetic geometry and one evaluation
1396 seed, and they do not include confidence intervals. The formal claims in the paper are the population
1397 statements and counterexamples proved in the main text and preceding appendices.

1398 **N LLM Usage**

1399 Large language models were used for editing, wording, and non-authoritative review-style feedback
1400 on presentation and claim calibration. All theoretical claims, proofs, experiments, and final text were
1401 checked and are the responsibility of the authors.

1402 **O Claim ledger**

1403 Table 4 summarizes results and perspective.

Table 4: Claim ledger and relation to prior work. The formal results concern population one-slot or residual-greedy sparse reconstruction, transfer to one-row SAE realizability, and the separation between decoder directions and latent support semantics.

Topic	What we prove	What we interpret	Connection to prior work
Population objective	For a fixed residual law and unit direction v , the optimal nonnegative one-slot code is $\langle (v, R) - \lambda \rangle_+$, with gain $U_\lambda(v; R) = \frac{1}{2} \mathbb{E}[\langle (v, R) - \lambda \rangle_+^2]$. For positive latent-ray residuals, this gain decomposes exactly over residual support events.	The one-slot sparse-reconstruction objective scores thresholded residual mass, not teacher-feature identity directly. Singleton, mixed, and broader support events contribute separately.	A population, free-direction analogue of thresholded matching pursuit, complementary to sparse recovery and dictionary-learning guarantees under incoherence, separability, or sparse-use assumptions. [22, 30, 14, 27, 2]
Teacher recovery	If the utility lost by rotating away from a singleton rank-one teacher event exceeds a conservative off-singleton competition budget, then all one-slot global maximizers lie near the residualized teacher direction.	This certifies recovery when clean singleton events dominate off-singleton contamination above threshold. Failure of the certificate is not claimed to imply failure of recovery.	Explains why restrictive sparsity, frequency, geometry, or hyperparameter alignment assumptions in SAE recovery theory help: they make teacher-identifying events dominate thresholded utility. [9, 6, 11, 7]
Merging and hedging	We give local and explicit population certificates under which rotating from a teacher direction toward another residual direction increases one-slot utility. In two-feature co-occurrence models, a mixed direction can beat individual teachers once joint events cross threshold.	Merging can arise from thresholded mixed-event mass before finite samples, amortized encoders, or optimization enter. Correlation is one source; the formal driver is above-threshold transverse residual mass.	Feature hedging and non-canonical-feature work show that trained SAEs can learn mixed latents. We isolate an oracle-level eventwise mechanism that rewards such directions. [5, 20]
Splitting	When a broad event family has heterogeneous residual directions, event-partition certificates show that a multi-direction specialized representation can obtain larger residual-greedy utility than forcing the family through one shared direction. We also separate this from encoder-induced splitting.	Splitting can reflect heterogeneous residual evidence, not only semantic oversegmentation or training instability. A semantic feature may arrive through several residual rays.	Connects empirical splitting and matching-pursuit-style residual decompositions to a population utility comparison. [8?, 10]
Absorption and residual paths	Sequential residual subtraction can suppress above-threshold evidence for a later teacher direction on part of its support. We also construct exact tied population residual paths where early choices lead to inequivalent later choices.	Absorption-like holes and residual-path alternatives can be substrates of the population residual-greedy objective. We do not claim a convergence theorem for SGD or a full account of empirical seed sensitivity.	Complements absorption, seed-instability, and SDL non-identifiability studies by isolating residual suppression and tie-breaking mechanisms. [8, 26, 28]
Encoder realizability	For a fixed decoder direction and leave-one-out residual, one encoder row performs squared-error regression onto the oracle target $d^\top R - \lambda$ with a nonnegative affine-ReLU output. The one-row value matches the oracle iff the oracle gate lies in the L^2 -closure of affine-ReLU gates; exact realization is the attained case.	Oracle preference and SAE realizability are separate layers. A direction may be oracle-preferred while a one-ReLU row cannot realize its gate exactly.	Identifies a representability gap between unrestricted sparse-coding gates and one-row ReLU SAE encoders. [3, 12, 9, 11]
Direction vs. support	Decoder direction alignment is a gated residual-moment condition. Under residualized linear independence, exact alignment with a teacher ray requires zero gated co-activation with all other residualized atoms. Counterexamples show that direction and support recovery do not imply each other.	A recovered decoder direction alone does not certify that the latent fires on the intended support. Directional semantics and activation-pattern semantics are distinct.	Refines latent-recovery work showing that reconstruction or decoder recovery need not identify latent variables. [28, 31]
Sparsity, frequency, diagnostics	Exact support recovery implies sample-wise sparsity and expected sparsity matching, but converses fail. Without latent supports, we prove observer-only residual-cap certificates and define direction, precision, recall, ambiguity, and gate-consistency diagnostics.	Sparsity and firing frequency are useful shadows of support recovery, but the objective depends on joint support events, amplitudes, thresholds, and geometry. Observer diagnostics certify residual geometry, not latent support semantics.	Explains what sparsity/frequency alignment captures in toy recovery studies and what it misses; gives population diagnostics for settings without ground-truth supports. [6, 9, 11, 3, 29]

1404 **NeurIPS Paper Checklist**

1405 **1. Claims**

1406 Question: Do the main claims made in the abstract and introduction accurately reflect the
1407 paper’s contributions and scope?

1408 Answer: [Yes]

1409 Justification: The abstract and introduction state theoretical claims about the population
1410 sparse-reconstruction objective, support-event decomposition, dominance regimes for re-
1411 covery, one-ReLU realizability, and direction/support separation. The paper does not claim
1412 empirical validation or a full theory of trained SAE optimization dynamics.

1413 Guidelines:

- 1414 • The answer [N/A] means that the abstract and introduction do not include the claims
1415 made in the paper.
- 1416 • The abstract and/or introduction should clearly state the claims made, including the
1417 contributions made in the paper and important assumptions and limitations. A [No] or
1418 [N/A] answer to this question will not be perceived well by the reviewers.
- 1419 • The claims made should match theoretical and experimental results, and reflect how
1420 much the results can be expected to generalize to other settings.
- 1421 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1422 are not attained by the paper.

1423 **2. Limitations**

1424 Question: Does the paper discuss the limitations of the work performed by the authors?

1425 Answer: [Yes]

1426 Justification: The paper discusses structural limitations including the positive latent-ray
1427 model, nonnegative coefficients, population rather than finite-sample analysis, fixed-residual
1428 or residual-greedy scope, non-vacuous residual charts, absence of an SGD convergence
1429 theory, lack of empirical validation on pretrained-model SAEs, and the limited illustrative
1430 role of the synthetic experiments.

1431 Guidelines:

- 1432 • The answer [N/A] means that the paper has no limitation while the answer [No] means
1433 that the paper has limitations, but those are not discussed in the paper.
- 1434 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 1435 • The paper should point out any strong assumptions and how robust the results are to
1436 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1437 model well-specification, asymptotic approximations only holding locally). The authors
1438 should reflect on how these assumptions might be violated in practice and what the
1439 implications would be.
- 1440 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1441 only tested on a few datasets or with a few runs. In general, empirical results often
1442 depend on implicit assumptions, which should be articulated.
- 1443 • The authors should reflect on the factors that influence the performance of the approach.
1444 For example, a facial recognition algorithm may perform poorly when image resolution
1445 is low or images are taken in low lighting. Or a speech-to-text system might not be
1446 used reliably to provide closed captions for online lectures because it fails to handle
1447 technical jargon.
- 1448 • The authors should discuss the computational efficiency of the proposed algorithms
1449 and how they scale with dataset size.
- 1450 • If applicable, the authors should discuss possible limitations of their approach to
1451 address problems of privacy and fairness.
- 1452 • While the authors might fear that complete honesty about limitations might be used by
1453 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
1454 limitations that aren’t acknowledged in the paper. The authors should use their best
1455 judgment and recognize that individual actions in favor of transparency play an impor-
1456 tant role in developing norms that preserve the integrity of the community. Reviewers
1457 will be specifically instructed to not penalize honesty concerning limitations.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509

3. Theory assumptions and proofs

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

Answer: [Yes]

Justification: The paper states the modeling assumptions, including finite-moment assumptions, the positive latent-ray model, nonnegative coefficients, residual-chart conditions, and one-ReLU encoder restrictions where relevant. The theoretical results are accompanied by theorem statements, separation constructions, and complete proofs in the appendix.

Guidelines:

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

4. Experimental result reproducibility

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

Answer: [Yes]

Justification: The synthetic experiments are illustrative rather than central to the main theoretical claims. The appendix specifies the data-generating process, SAE architecture, objective, optimizer, training steps, validation protocol, L_1 selection procedure, seeds, evaluation set size, matching metrics, and aggregate results needed to reproduce the two-stage SAE experiments.

Guidelines:

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

- 1510 (c) If the contribution is a new model (e.g., a large language model), then there should
1511 either be a way to access this model for reproducing the results or a way to reproduce
1512 the model (e.g., with an open-source dataset or instructions for how to construct
1513 the dataset).
- 1514 (d) We recognize that reproducibility may be tricky in some cases, in which case
1515 authors are welcome to describe the particular way they provide for reproducibility.
1516 In the case of closed-source models, it may be that access to the model is limited in
1517 some way (e.g., to registered users), but it should be possible for other researchers
1518 to have some path to reproducing or verifying the results.

1519 5. Open access to data and code

1520 Question: Does the paper provide open access to the data and code, with sufficient instruc-
1521 tions to faithfully reproduce the main experimental results, as described in supplemental
1522 material?

1523 Answer: [No]

1524 Justification: The paper does not currently provide an anonymized code release. The
1525 experiments use synthetic data generated from fully specified distributions, and the appendix
1526 provides the information needed to reimplement them.

1527 Guidelines:

- 1528 • The answer [N/A] means that paper does not include experiments requiring code.
- 1529 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
1530 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1531 • While we encourage the release of code and data, we understand that this might not
1532 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
1533 including code, unless this is central to the contribution (e.g., for a new open-source
1534 benchmark).
- 1535 • The instructions should contain the exact command and environment needed to run to
1536 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
1537 //neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1538 • The authors should provide instructions on data access and preparation, including how
1539 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1540 • The authors should provide scripts to reproduce all experimental results for the new
1541 proposed method and baselines. If only a subset of experiments are reproducible, they
1542 should state which ones are omitted from the script and why.
- 1543 • At submission time, to preserve anonymity, the authors should release anonymized
1544 versions (if applicable).
- 1545 • Providing as much information as possible in supplemental material (appended to the
1546 paper) is recommended, but including URLs to data and code is permitted.

1547 6. Experimental setting/details

1548 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
1549 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1550 Answer: [Yes]

1551 Justification: The appendix gives the experimental setting: synthetic generator, dimensionality,
1552 number of features, support and magnitude distributions, evaluation sample size, SAE
1553 architecture, reconstruction and L_1 objective, Adam hyperparameters, training schedule,
1554 validation protocol, decoder normalization, L_1 binary search, seeds, and matching metrics.

1555 Guidelines:

- 1556 • The answer [N/A] means that the paper does not include experiments.
- 1557 • The experimental setting should be presented in the core of the paper to a level of detail
1558 that is necessary to appreciate the results and make sense of them.
- 1559 • The full details can be provided either with the code, in appendix, or as supplemental
1560 material.

1561 7. Experiment statistical significance

1562 Question: Does the paper report error bars suitably and correctly defined or other appropriate
1563 information about the statistical significance of the experiments?

1564 Answer: [No]

1565 Justification: The reported experiments are small synthetic illustrations and do not include
1566 formal confidence intervals or error bars.

1567 Guidelines:

- 1568 • The answer [N/A] means that the paper does not include experiments.
- 1569 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
1570 intervals, or statistical significance tests, at least for the experiments that support the
1571 main claims of the paper.
- 1572 • The factors of variability that the error bars are capturing should be clearly stated (for
1573 example, train/test split, initialization, random drawing of some parameter, or overall
1574 run with given experimental conditions).
- 1575 • The method for calculating the error bars should be explained (closed form formula,
1576 call to a library function, bootstrap, etc.)
- 1577 • The assumptions made should be given (e.g., Normally distributed errors).
- 1578 • It should be clear whether the error bar is the standard deviation or the standard error
1579 of the mean.
- 1580 • It is OK to report 1-sigma error bars, but one should state it. The authors should
1581 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
1582 of Normality of errors is not verified.
- 1583 • For asymmetric distributions, the authors should be careful not to show in tables or
1584 figures symmetric error bars that would yield results that are out of range (e.g., negative
1585 error rates).
- 1586 • If error bars are reported in tables or plots, the authors should explain in the text how
1587 they were calculated and reference the corresponding figures or tables in the text.

1588 8. Experiments compute resources

1589 Question: For each experiment, does the paper provide sufficient information on the com-
1590 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1591 the experiments?

1592 Answer: [No]

1593 Justification: The experiments are lightweight synthetic SAE runs, but the current manuscript
1594 does not fully specify wall-clock time, hardware type, memory, or total compute. Since the
1595 experiments are illustrative and not the main contribution, this omission does not affect the
1596 reproducibility of the theoretical results.

1597 Guidelines:

- 1598 • The answer [N/A] means that the paper does not include experiments.
- 1599 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1600 or cloud provider, including relevant memory and storage.
- 1601 • The paper should provide the amount of compute required for each of the individual
1602 experimental runs as well as estimate the total compute.
- 1603 • The paper should disclose whether the full research project required more compute
1604 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1605 didn't make it into the paper).

1606 9. Code of ethics

1607 Question: Does the research conducted in the paper conform, in every respect, with the
1608 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1609 Answer: [Yes]

1610 Justification: The work is theoretical and uses only synthetic toy experiments. It does not
1611 involve human subjects, private data, scraped datasets, deployed systems, or release of
1612 high-risk models.

1613 Guidelines:

- 1614
- 1615
- 1616
- 1617
- 1618
- 1619
- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

1620

1621

1622

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

1623

Answer: [Yes]

1624

1625

1626

1627

Justification: The paper is foundational work on interpretability and sparse representations. Its positive impact is improved understanding of when SAE features should or should not be interpreted as ground-truth features; a potential negative impact is misplaced confidence in interpretability tools if their structural limitations are ignored.

1628

Guidelines:

- 1629
- 1630
- 1631
- 1632
- 1633
- 1634
- 1635
- 1636
- 1637
- 1638
- 1639
- 1640
- 1641
- 1642
- 1643
- 1644
- 1645
- 1646
- 1647
- 1648
- 1649
- 1650
- The answer [N/A] means that there is no societal impact of the work performed.
 - If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

1651

1652

1653

1654

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

1655

Answer: [No]

1656

1657

Justification: The paper does not release pretrained models, image generators, scraped datasets, or other high-risk assets.

1658

Guidelines:

- 1659
- 1660
- 1661
- 1662
- 1663
- 1664
- 1665
- The answer [N/A] means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.

- 1666 • We recognize that providing effective safeguards is challenging, and many papers do
1667 not require this, but we encourage authors to take this into account and make a best
1668 faith effort.

1669 12. Licenses for existing assets

1670 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1671 the paper, properly credited and are the license and terms of use explicitly mentioned and
1672 properly respected?

1673 Answer: [N/A]

1674 Justification: The paper does not rely on external datasets, pretrained models, or proprietary
1675 assets. Existing scholarly work is credited through citations.

1676 Guidelines:

- 1677 • The answer [N/A] means that the paper does not use existing assets.
- 1678 • The authors should cite the original paper that produced the code package or dataset.
- 1679 • The authors should state which version of the asset is used and, if possible, include a
1680 URL.
- 1681 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1682 • For scraped data from a particular source (e.g., website), the copyright and terms of
1683 service of that source should be provided.
- 1684 • If assets are released, the license, copyright information, and terms of use in the
1685 package should be provided. For popular datasets, paperswithcode.com/datasets
1686 has curated licenses for some datasets. Their licensing guide can help determine the
1687 license of a dataset.
- 1688 • For existing datasets that are re-packaged, both the original license and the license of
1689 the derived asset (if it has changed) should be provided.
- 1690 • If this information is not available online, the authors are encouraged to reach out to
1691 the asset's creators.

1692 13. New assets

1693 Question: Are new assets introduced in the paper well documented and is the documentation
1694 provided alongside the assets?

1695 Answer: [N/A]

1696 Justification: The paper does not introduce a new dataset, benchmark, pretrained model, or
1697 reusable software package as a contribution.

1698 Guidelines:

- 1699 • The answer [N/A] means that the paper does not release new assets.
- 1700 • Researchers should communicate the details of the dataset/code/model as part of their
1701 submissions via structured templates. This includes details about training, license,
1702 limitations, etc.
- 1703 • The paper should discuss whether and how consent was obtained from people whose
1704 asset is used.
- 1705 • At submission time, remember to anonymize your assets (if applicable). You can either
1706 create an anonymized URL or include an anonymized zip file.

1707 14. Crowdsourcing and research with human subjects

1708 Question: For crowdsourcing experiments and research with human subjects, does the paper
1709 include the full text of instructions given to participants and screenshots, if applicable, as
1710 well as details about compensation (if any)?

1711 Answer: [N/A]

1712 Justification: The paper does not involve crowdsourcing, human-subject studies, user studies,
1713 or participant compensation.

1714 Guidelines:

- 1715 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
1716 with human subjects.

1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755

- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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

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

Answer: [N/A]

Justification: The paper does not involve human subjects, crowdsourcing, or participant data, so IRB or equivalent approval is not applicable.

Guidelines:

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

16. Declaration of LLM usage

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

Answer: [Yes]

Justification: The paper declares LLM usage for editing, wording, and non-authoritative review-style feedback.

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.