
Crosscoding Through Time: Sparse Feature Discovery Across Sequence Positions

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Dictionary learning methods - such as Sparse Autoencoders (SAEs) and cross-
2 coders - decompose model activations into human-interpretable building blocks.
3 We introduce *temporal crosscoders*, a simple and flexible framework for feature
4 discovery in Large Language Models (LLMs). To properly evaluate temporal
5 crosscoders we develop TempBench: a panel of synthetic and real-world tasks for
6 evaluating temporal structures. Temporal crosscoders outperform both conventional
7 and temporal architectures in both of our synthetic settings and on two out of four
8 of the real world settings - ahead of other candidate architectures. Most strikingly,
9 they can detect backtracking - a key reasoning behavior - at a 40% higher rate than
10 conventional SAEs, and are 15% more effective in inducing it. Our results establish
11 temporal crosscoders as a simple and flexible framework for feature discovery, both
12 local and temporal. We provide full code at the following anonymous repository:
13 <https://anonymous.4open.science/r/temp-bench-anon/>.

14 1 Introduction

15 Sparse autoencoders have become a foundational tool for interpreting large language models. By
16 decomposing activations into sparse, reusable directions, SAEs provide a feature dictionary that
17 supports monosemantic feature discovery, sparse probing, steering, model comparison, and large-scale
18 interpretability workflows [Cunningham et al., 2023, Bricken et al., 2023a].

19 However, standard SAEs reconstruct each token independently of nearby positions. This assumes
20 that the relevant feature structure is local in sequence position [Lubana et al., 2025]. Natural language
21 does not satisfy this assumption. It has hierarchical syntactic structure [Chomsky, 1956]; tokens
22 are constrained by expectations from prior context [Levy, 2008]; and neural responses to language
23 integrate information over multiple temporal receptive windows, from short phrases to longer narrative
24 contexts [Lerner et al., 2011].

25 Recent work has made this mismatch explicit, framing it as an architectural issue: standard SAE
26 objectives factor over positions and so impose an independence prior over time, missing latent
27 structure shared across tokens [Lubana et al., 2025, Bhalla et al., 2026]. A complementary observation
28 is that correlations between sequence positions decay approximately as a power law with distance
29 [Cagnetta and Wyart, 2024], so the relevant dependencies are not confined to a single fixed scale.
30 Feature discovery therefore needs a windowed architecture that can capture correlations within a
31 local context while allowing learned features to live at different scales inside that window.

32 In this work, we make four contributions to temporal feature discovery:

- 33 1. We introduce temporal crosscoders section 3, a simple and flexible architecture for capturing
34 temporal structure in LLMs.

- 35 2. We generalize existing synthetic benchmarks for ground-truth feature recovery to the tempo-
36 ral setting section 4.
- 37 3. We assemble four real-world evaluations with clear behavioral signals for comparing
38 dictionary-learning architectures section 5.
- 39 4. We combine these synthetic and real-world evaluations into **TempBench**, an evaluation
40 suite for conventional and temporal SAE architectures, and use it to benchmark the current
41 architectural frontier under matched conditions.

42 Across *TempBench*, temporal aggregation helps when the target feature is supported by evidence
43 distributed over nearby tokens, and hurts on tasks where per-token or length-correlated signals
44 dominate.

45 2 Related Work

46 **Sparse autoencoders and synthetic benchmarks.** Sparse autoencoders (SAEs) trained on the
47 residual stream have become the standard tool for decomposing language-model activations into
48 a sparse basis of features [Cunningham et al., 2023, Bricken et al., 2023a, Templeton et al., 2024,
49 Gao et al., 2024, Bussmann et al., 2024, Lieberum et al., 2024]. Because SAE features are derived
50 rather than ground-truth, evaluating *which* dictionary recovers *what* requires either downstream
51 proxies (loss recovery, sparse probing [Kantamneni et al., 2025, Bhalla et al., 2026]) or synthetic data
52 with known features. Chanin [2025] establish the canonical synthetic setting: ground-truth feature
53 directions $\{\mathbf{f}_i\}$ are fired according to a distribution that factorizes across sequence positions, so the
54 recoverable ground truth is purely local. This setup has been used to study superposition and feature
55 absorption [Makelov et al., 2024, Karvonen et al., 2024], but by construction it cannot test whether a
56 dictionary recovers structure that is shared across positions, because no such structure exists in the
57 data. Our section 4 generalizes this setting by promoting the firing distribution to a Hidden Markov
58 Model, which (a) preserves the local synthetic setting as a special case and (b) admits two distinct
59 ground truths — emission directions and hidden-state directions — that dissociate per-position from
60 cross-position recovery.

61 **Crosscoders.** Crosscoders [Lindsey et al., 2024] extend SAEs by encoding and decoding across an
62 additional axis: a single sparse latent vector explains multiple inputs simultaneously, with axis-specific
63 decoder weights. The original proposal targets the layer axis (a single feature reconstructs activations
64 at several layers) and motivates the architecture by model-diffing applications — comparing base and
65 finetuned models, or models from different training runs [Minder et al., 2026, Lindsey et al., 2024] —
66 and cross-layer feature discovery. To our knowledge, the crosscoder framework has not been applied
67 to the sequence axis. Our temporal crosscoder (section 3) is the temporal analogue.

68 **Temporal sparse dictionaries.** Two recent proposals share our motivation that per-token dictionar-
69 ies fragment multi-token concepts, and we use both as baselines. The temporal SAE of Bhalla et al.
70 [2025] introduces a T -position window with a per-position encoder/decoder and a multi-distance
71 contrastive loss that aligns features describing the same temporal pattern across nearby offsets. Per-
72 position weights mean that T-SAE features are still anchored to single positions; the contrastive loss
73 is what couples them across the window. The temporal feature analyzer (TFA) of Lubana et al. [2025]
74 takes a different route, introducing a temporal prior on the dictionary itself. The temporal crosscoder
75 differs from both by tying every position of the window to a single shared latent and letting full-rank
76 cross-position decoders distribute that latent’s mass.

77 **Systematic benchmarking.** Each of the temporal architectures above was introduced and evaluated
78 in isolation. There is no existing study that evaluates the candidate architectures side-by-side on a
79 panel of comparable tasks. The closest prior efforts are SAEBench [Karvonen et al., 2025] and the
80 sparse-probing panel of Kantamneni et al. [2025], both of which standardise downstream evaluation
81 but cover only per-token dictionaries.

82 3 Temporal Crosscoders

83 Let $\mathbf{x}_t \in \mathbb{R}^d$ denote the residual-stream activation at position t , and let $\mathbf{X}_t = (\mathbf{x}_t, \dots, \mathbf{x}_{t+T-1})$ be a
84 temporal window. A temporal crosscoder learns one sparse latent vector $\mathbf{u}_t \in \mathbb{R}^H$ shared across the

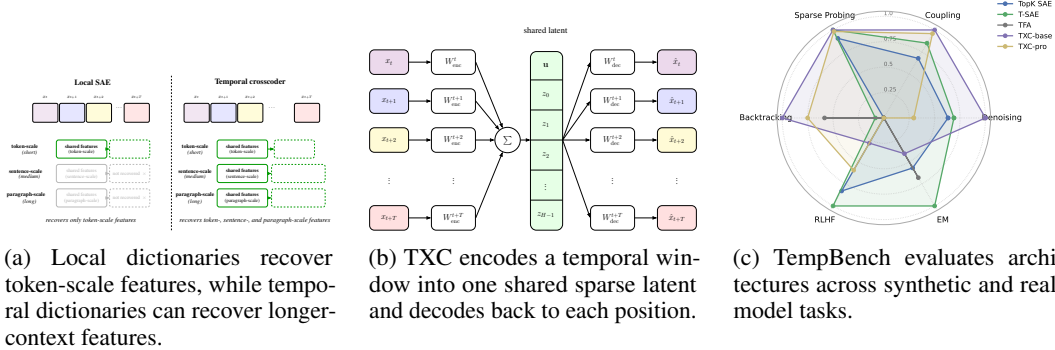


Figure 1: Schematic motivation, architecture, and evaluation suite for temporal crosscoders.

85 whole window:

$$\mathbf{u}_t = \sigma \left(\sum_{\tau=0}^{T-1} W_{\text{enc}}^{(\tau)} \mathbf{x}_{t+\tau} + \mathbf{b}_{\text{enc}} \right), \quad \hat{\mathbf{x}}_{t+\tau} = W_{\text{dec}}^{(\tau)} \mathbf{u}_t + \mathbf{b}_{\text{dec}}^{(\tau)}.$$

86 This is the temporal analogue of a layerwise acausal crosscoder [Lindsey et al., 2024]: the layer index
 87 is replaced by a relative position index. A feature may reconstruct primarily one token, recovering
 88 local-SAE behavior, or distribute decoder mass across several positions, representing a temporally
 89 extended pattern.

90 We train temporal crosscoders with a reconstruction loss over the full window, together with explicit
 91 sparsity from the latent activation function. In our default implementation, σ is BatchTopK [Bussmann
 92 et al., 2024], a variant of TopK [Gao et al., 2024], applied to the shared latent vector for each window,
 93 and we set decoder biases to zero so that reconstruction mass is assigned to features rather than biases.
 94 We define and use two TXC variants in this work; definitions and model decisions are explained in
 95 section A. Our base temporal crosscoder TXC-base uses only this acausal window-reconstruction
 96 objective. The larger TXC-pro variant adds two temporal inductive biases: a Matryoshka objective
 97 over nested latent groups, encouraging useful feature dictionaries at multiple effective widths, and a
 98 multi-distance contrastive loss over shifted windows, encouraging features that describe the same
 99 temporal pattern to align across nearby offsets.

100 3.1 Compared Architectures and Matching

101 Across TempBench, we compare TXC-base and TXC-pro against three baselines: TopK SAE [Gao
 102 et al., 2024], T-SAE [Bhalla et al., 2025], and TFA [Lubana et al., 2025]. These cover local dictionar-
 103 ies, temporal dictionaries with position-local structure, temporal-prior methods, and crosscoder-style
 104 sharing along a non-temporal axis.

105 Unless otherwise stated, all architectures are trained on the same activation cache and compared at
 106 fixed hookpoint, dictionary width, and expected sparsity. We write k_{pos} for the active-latent budget
 107 per sequence position: a per-token dictionary uses k_{pos} active latents at each token, while a windowed
 108 architecture of length T receives the matched window budget Tk_{pos} .

109 4 Benchmarking Temporal Crosscoders in a Synthetic Setting

110 Synthetic datasets make it possible to evaluate feature recovery against known ground-truth features
 111 [Chanin, 2025]. Existing synthetic benchmarks define feature activations independently at each
 112 sequence position, so the ground-truth feature structure is local. To benchmark temporal architectures,
 113 we generalize this setting by making the feature-firing pattern a stochastic process across sequence
 114 positions.

115 Consider a synthetic activation sequence

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T), \quad \mathbf{x}_t \in \mathbb{R}^d,$$

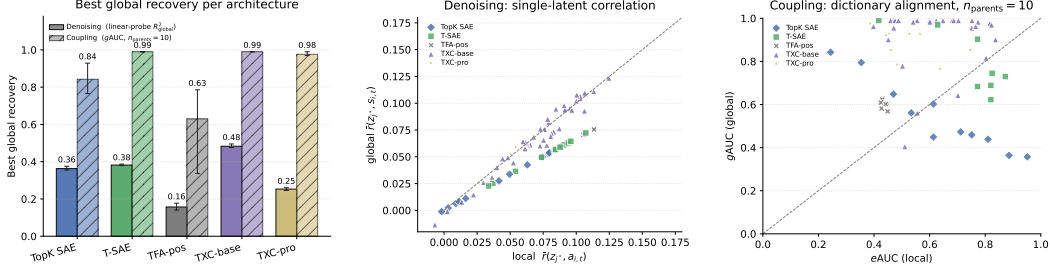


Figure 2: Synthetic global-feature recovery. *Left*: best global recovery per architecture, using linear-probe R_{global}^2 for Denoising and $gAUC$ for Coupling. *Middle*: Denoising single-latent correlations with noisy emissions $a_{i,t}$ versus clean hidden states $s_{i,t}$; the diagonal marks equal local and global correlation. *Right*: Coupling dictionary alignment with emission directions ($eAUC$) versus hidden-chain footprints ($gAUC$). Temporal architectures recover more global structure in Coupling, while TXC-base gives the strongest Denoising recovery.

116 defined over a temporal window of length T . A synthetic setting is specified by three pieces of data.
 117 First, a set of N ground-truth feature directions

$$\mathcal{F} = \{\mathbf{f}_1, \dots, \mathbf{f}_N\}, \quad \mathbf{f}_i \in \mathbb{R}^d.$$

118 Second, a probability distribution P_T over binary feature-firing masks

$$A = (\mathbf{a}_1, \dots, \mathbf{a}_T) \in (\{0, 1\}^N)^T,$$

119 where

$$\mathbf{a}_t = (a_{t,1}, \dots, a_{t,N})$$

120 and $a_{t,i} = 1$ indicates that feature i fires at position t . Third, a distribution M_T over feature
 121 magnitudes

$$M = (\mathbf{m}_1, \dots, \mathbf{m}_T), \quad \mathbf{m}_t = (m_{t,1}, \dots, m_{t,N}).$$

122 A synthetic activation is then constructed as

$$\mathbf{x}_t = \sum_{i=1}^N a_{t,i} m_{t,i} \mathbf{f}_i.$$

123 Thus, the feature directions specify what can be recovered, while P_T specifies the local and temporal
 124 structure of when those features appear. The conventional synthetic setting is the special case where
 125 P_T and M_T factorize across sequence positions. In our setting, local correlations are encoded by the
 126 one-position distribution over \mathbf{a}_t , while temporal correlations are encoded by the joint distribution
 127 over $\mathbf{a}_{1:T}$.

128 A natural way to specify the firing-mask distribution P_T is by an edge-emitting Hidden Markov
 129 Model (HMM), following the convention used in computational mechanics [Crutchfield and Young,
 130 1989, Crutchfield, 2012, Shai et al., 2024, 2026]. An HMM \mathcal{H} consists of

$$\mathcal{H} = (\mathcal{X}, \mathcal{S}, \boldsymbol{\eta}^\theta, \{T^{(x)}\}_{x \in \mathcal{X}}),$$

131 where the observed alphabet is the set of local firing masks,

$$\mathcal{X} = \{0, 1\}^N.$$

132 Thus, each emitted symbol $x \in \mathcal{X}$ is a full binary vector indicating which features fire at a single
 133 sequence position. The set \mathcal{S} is the latent state space, $\boldsymbol{\eta}^\theta \in \Delta(\mathcal{S})$ is the initial distribution over latent
 134 states, and $T^{(x)}$ is the symbol-labeled transition matrix with entries

$$T_{h,h'}^{(x)} = \Pr(h_{t+1} = h', \mathbf{a}_t = x \mid h_t = h).$$

135 The HMM generates the distribution over firing patterns by

$$P_T(\mathbf{a}_{1:T}) = \boldsymbol{\eta}^{\theta \top} T^{(\mathbf{a}_1)} T^{(\mathbf{a}_2)} \dots T^{(\mathbf{a}_T)} \mathbf{1}.$$

136 A dictionary learning architecture can be evaluated at two levels. At the local level, we ask whether
 137 its decoder directions recover the emitted feature directions \mathbf{f}_i and whether its activations track the
 138 local firing indicators $a_{t,i}$. At the global level, we ask whether its learned latents track the HMM state,
 139 or the corresponding hidden feature directions induced by those states. In non-synthetic settings, this
 140 corresponds to the distinction between learning local cues associated with a behavior and learning a
 141 latent feature for the behavior itself.

142 We benchmark dictionary learning architectures on two HMM special cases in which the observed
 143 local feature firings are only indirect evidence for a cleaner hidden temporal process. The two benches
 144 differ in where the local/global distinction appears. In the **Denoising** bench, each emitted feature has
 145 its own hidden Markov chain, so there is a one-to-one correspondence between hidden states $s_{i,t}$ and
 146 noisy local firings $a_{i,t}$. Since both variables are associated with the same residual-stream direction
 147 \mathbf{f}_i , local-vs-global recovery cannot be read off from decoder alignment; it must be read off from
 148 the learned activation traces. In the **Coupling** bench, a smaller set of hidden chains drives a larger
 149 emission alphabet through OR-gate coupling and per-token noise. There, local emission directions
 150 and hidden-chain footprints are geometrically distinct, so local-vs-global recovery can be measured
 151 directly from decoder alignment.

152 **The Denoising bench.** The Denoising bench is a panel of $N = 20$ independent two-state Markov
 153 chains, one per feature. Each chain has a sticky hidden state $s_{i,t}$ with self-correlation $\rho = 0.7$, and
 154 emits a noisy firing $a_{i,t}$ through a Bernoulli channel with hit rate $p_B = 0.625$; the activation reads
 155 out the noisy firings,

$$\mathbf{x}_t = \sum_{i=1}^N a_{i,t} m_{i,t} \mathbf{f}_i.$$

156 Thus the clean hidden state $s_{i,t}$ and noisy emission $a_{i,t}$ differ in their time traces, not in their per-
 157 position residual-stream direction. A learned feature whose decoder points along \mathbf{f}_i may either track
 158 the noisy emission $a_{i,t}$ or use temporal context to recover the cleaner hidden state $s_{i,t}$.

159 We evaluate this distinction at the level of the matched feature’s activation trace. For each ground-truth
 160 feature i , let $j(i)$ be the learned feature whose decoder slice has the highest cosine similarity with \mathbf{f}_i .
 161 We define

$$r_{\text{local},i} = \text{CORR}(z_{j(i)}, a_i), \quad r_{\text{global},i} = \text{CORR}(z_{j(i)}, s_i),$$

162 where correlations are taken over held-out sequence positions. The single-latent scatter in section 4
 163 plots $(r_{\text{local},i}, r_{\text{global},i})$. A per-token model that merely tracks the noisy emission should lie near the
 164 no-denoising floor

$$r_{\text{global},i} \approx \text{CORR}(a_i, s_i) r_{\text{local},i}.$$

165 The diagonal $y = x$ is a stronger reference line: points above it are more correlated with the clean
 166 hidden state than with the noisy emission. We also report a full-code linear probe, ridge-regressing
 167 $s_{i,t}$ on the learned latent code and reporting R_{global}^2 in the headline bar chart. In section 4, per-token
 168 baselines remain near the no-denoising floor. TXC-base moves above this floor as the window length
 169 grows and achieves the best full-code denoising score ($R_{\text{global}}^2 = 0.48$). TXC-pro does not improve
 170 on this setting, suggesting that the additional temporal objectives are not uniformly helpful for this
 171 simple hidden-state recovery task.

172 **The Coupling bench.** The Coupling bench removes the one-to-one correspondence between hidden
 173 chains and emitted features. A set of $K = 10$ hidden Markov chains drives a larger alphabet of
 174 $M = 20$ local emissions through a binary coupling matrix $C \in \{0, 1\}^{M \times K}$. At each position,
 175 emission j is active when at least one of its parent chains is active, followed by a Bernoulli noise
 176 channel:

$$a_{j,t} \sim \text{Bernoulli} \left(p_B \mathbb{1} \left[\sum_{k=1}^K C_{jk} h_{k,t} \geq 1 \right] \right),$$

177 with $p_A = 0$ in the off state. The observed activation is

$$\mathbf{x}_t = \sum_{j=1}^M a_{j,t} m_{j,t} \mathbf{f}_j^{\text{emit}}.$$

178 This setting has two spatial ground truths. The local ground truth is the set of emission directions

$$\{\mathbf{f}_j^{\text{emit}}\}_{j=1}^M,$$

179 which are directly written into the activation. The global ground truth is the set of hidden-chain
180 footprints

$$\mathbf{f}_k^{\text{hid}} = \text{normalize} \left(\sum_{j:C_{jk}=1} \mathbf{f}_j^{\text{emit}} \right), \quad k = 1, \dots, K.$$

181 Unlike in the Denoising bench, local and global recovery are now visible in the dictionary geometry
182 itself: a locally aligned dictionary points at emission directions, while a globally aligned dictionary
183 points at hidden-chain footprints. We therefore evaluate the Coupling bench with decoder-alignment
184 AUCs. The local score $e\text{AUC}$ is the feature-recovery AUC against the M emission directions, and
185 the global score $g\text{AUC}$ is the feature-recovery AUC against the hidden-chain footprints.

186 We study coupling-density regimes that vary how much the hidden-chain footprints overlap. The
187 moderate-overlap regime takes $n_{\text{parents}} = 5$, so each emission is driven by a random half of the
188 hidden chains while different hidden chains still have distinct child sets. We also include a dense-
189 overlap stress test with $n_{\text{parents}} = 10$. Since this couples every emission to every hidden chain, the
190 individual hidden chain footprints collapse to a shared aggregate direction; we use this regime as a
191 local-vs-global stress test rather than as a disentanglement benchmark. Both regimes fix $K = 10$,
192 $M = 20$, $\rho = 0.9$, $\pi = 0.05$, $p_B = 0.5$, $d = 256$, $d_{\text{sae}} = 40$, and sweep T and k_{pos} over three seeds.
193 The architectures are the same set used on the Denoising bench.

194 **Coupling bench results.** In the dense-overlap stress test, temporal architectures recover the aggregate
195 hidden footprint while per-token-only architectures are pulled toward local emission directions.
196 T-SAE and TXC-base both reach $g\text{AUC} = 0.99$ at $k_{\text{pos}} = 1$, and TXC-pro reaches 0.98. The TopK
197 SAE follows the opposite trend, dropping from $g\text{AUC} = 0.92$ at $k_{\text{pos}} = 1$ to 0.44 at $k_{\text{pos}} = 8$ as it
198 reallocates dictionary capacity from the aggregate hidden footprint to the many local emission direc-
199 tions. The decoder-alignment scatter in section 4 makes this tradeoff concrete: temporal architectures
200 remain high in global alignment across a wide range of local-alignment scores, while TopK SAE
201 traces a down-right path as k_{pos} grows. Across the full $(p_B, n_{\text{parents}})$ sweep, the temporal advantage
202 is largest at moderate-to-high coupling overlap, with a peak $g\text{AUC}$ gap of ≈ 0.47 .

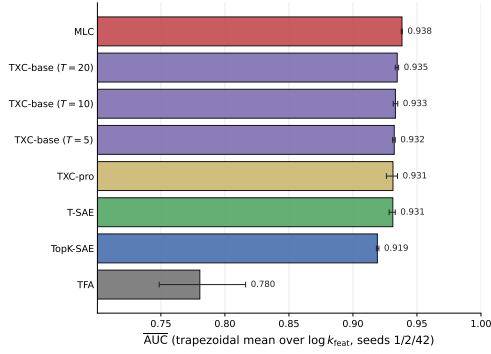
203 5 Real-Model Evaluations

204 The synthetic benchmarks above test whether a dictionary can recover known temporal ground truth.
205 We now ask whether the same temporal inductive bias is useful in trained language models, where
206 the relevant features are not known in advance and must be judged through downstream behavior. We
207 evaluate four complementary settings: sparse probing for static concept readout, backtracking for
208 causal control of a multi-token reasoning behavior, emergent misalignment for steering and detecting
209 a finetuned behavioral shift, and HH-RLHF preference decomposition as a negative control where
210 response length is a strong confounder. Across these settings, temporal aggregation helps when
211 the target behavior is distributed across nearby tokens, but can hurt when the benchmark rewards
212 length-correlated or per-token features. We include qualitative autointerpretation and per-token
213 activation visualizations in section G.

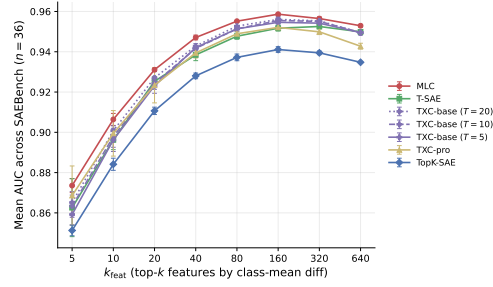
214 5.1 Sparse Probing

215 Sparse probing tests whether dictionary features support simple classifiers for human-interpretable
216 text concepts. We evaluate `google/gemma-2-2b-it` layer 13 on the 36-task SAEBench panel, with
217 all architectures trained at matched expected per-token $L_0 = 20$. For each task and feature budget
218 $k_{\text{feats}} \in \{5, 10, 20, 40, 80, 160, 320, 640\}$, we select features by training-split class-mean difference
219 and fit an ℓ_1 -regularized logistic probe. Full architecture, pooling, task, and aggregation details are in
220 section C.

221 Figure 3 shows that MLC leads at $\overline{\text{AUC}} = 0.907$, followed by TXC-base at 0.899–0.902, T-
222 SAE/TXC-pro at 0.897–0.899, and TopK SAE at 0.886. The temporal-family gap over TopK
223 SAE is small but larger than seed variation, appears by $k_{\text{feats}} \geq 20$, and saturates near 160. We read
224 this as evidence that sparse probing benefits from temporal-window features, but that most of the
225 gain comes from having temporal aggregation at all rather than from a specific TXC variant.

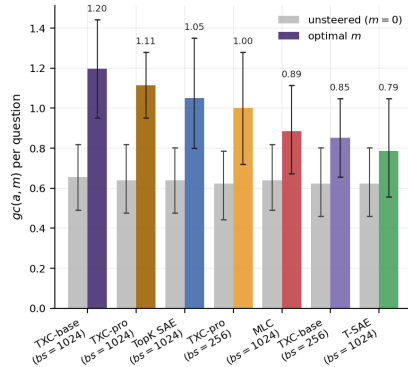


(a) $\overline{\text{AUC}}$ across the feature-budget sweep.

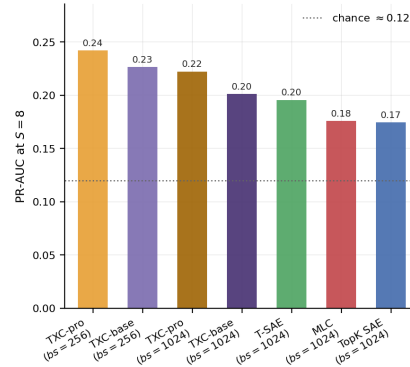


(b) AUC vs. number of probe features, with TFA omitted for scale.

Figure 3: Sparse probing on `google/gemma-2-2b-it` layer 13 over 36 SAEBench tasks. Temporal methods cluster near the top in $\overline{\text{AUC}}$ and outperform the per-token TopK SAE baseline; the advantage appears by $k_{\text{feats}} \geq 20$ and saturates near 160. Full curves including TFA are in fig. 7.



(a) Inducement: genuine backtracking at each architecture’s optimal steering magnitude.



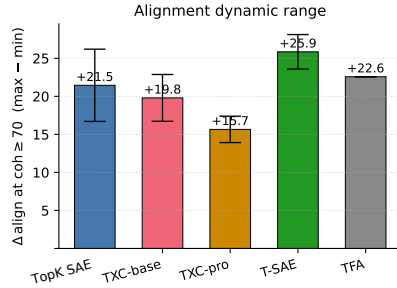
(b) Detection: sparse-probe PR-AUC at top- $S = 8$ features.

Figure 4: Backtracking inducement and detection on DeepSeek-R1-Distill-Llama-8B. TXC-base gives the strongest causal inducement, while TXC-pro gives the best sparse-probe detection above the class-prior baseline. A qualitative example is shown in fig. 12.

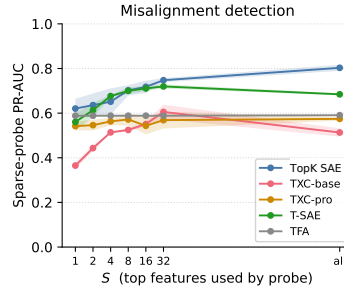
226 5.2 Inducing and Detecting Backtracking in a Reasoning Model

227 Reasoning models exhibit backtracking: the model commits to a partial argument, notices a problem,
 228 and reverses course. Because this signature spans multiple tokens, it is a natural testbed for temporal-
 229 window features. We follow the steering-vector setup of Ward et al. [2025], mining directions from
 230 Llama-3.1-8B and applying them at layer 10 of DeepSeek-R1-Distill-Llama-8B. The cohort
 231 contains 61 MATH-500 questions, and the cut25 protocol resamples continuations under steering
 232 across a grid of magnitudes.

233 We evaluate both causality and detection. Inducement counts genuine backtracking events under
 234 a Sonnet-4.6 judge, reported against the shared unsteered floor; detection trains sparse probes
 235 to distinguish backtracking from non-backtracking continuations under grouped cross-validation.
 236 Figure 4 shows that TXC-base gives the strongest causal inducement, while TXC-pro gives the
 237 best sparse-probe detection above the class-prior baseline. Full magnitude curves, judge rubric, and
 238 robustness checks are in section F.



(a) Steering dynamic range at coh ≥ 70 .



(b) Sparse-probe PR-AUC at $S = 16$.

Figure 5: Steering and detection for emergent misalignment in Qwen-2.5-7B-Instruct with a bad-medical-advice LoRA. T-SAE performs best on both axes, while TXC variants underperform; full α -frontiers are in figs. 8 and 9.

239 5.3 Emergent Misalignment

240 Emergent misalignment (EM) occurs when a model finetuned on a narrow misaligned task generalizes
 241 to broader misaligned behavior. Since EM responses are coherent multi-token behaviors and have
 242 been controlled by steering on residual-stream directions, we compare architectures on two axes:
 243 steering control and sparse-probe detection.

244 We evaluate Qwen-2.5-7B-Instruct finetuned with a bad-medical-advice LoRA. For steering, we
 245 adapt the Wang-stage procedure to select candidate features and report the maximum alignment
 246 dynamic range achievable above a coherence threshold of 70. For detection, we train sparse probes
 247 to distinguish aligned from misaligned rollouts and report PR-AUC. Figure 5 shows the opposite
 248 pattern from backtracking: T-SAE performs best on both axes, while TXC variants underperform.
 249 Full steering frontiers are in section D.

250 5.4 Decomposing HH-RLHF Preference Data

251 HH-RLHF pairs each prompt with a preferred chosen and dispreferred rejected response Bai et al.
 252 [2022]. Since Bhalla et al. [2025] use this dataset to test whether temporal features recover multi-
 253 token preference signals, we rerun their §4.5 / Appendix B.1 protocol on our locked architectures.
 254 We rank features by $\mu_{\text{rejected}} - \mu_{\text{chosen}}$ over response tokens and classify top features by correlation
 255 with response-length difference.

256 The result is a negative case for TXC. T-SAE at the paper setting $k = 20$ recovers 14/20 semantic top
 257 features, no length-spurious features, and 63% semantic mass. At the larger $k_{\text{win}} = 500$ budget, TopK
 258 SAE and T-SAE both fall to roughly 10/20 semantic features and 50% semantic mass. TXC performs
 259 worst on this static decomposition metric, with 7/20 semantic and 3 length-spurious features. We read
 260 this as evidence that temporal aggregation can help causal interventions spanning multiple tokens, as
 261 in section 5.2, while hurting rank-based interpretability when response length is a confounder. Full
 262 summary plots, scatter plots, and autointerpreted labels are in section E.

263 6 Limitations

264 Our work has two key limitations. Although we believe TempBench provides useful guidance to the
 265 dictionary-learning community, it is not a full map of how each architecture behaves. Our results,
 266 summarized in fig. 1c, tentatively suggest that different temporal architectures, namely TXCs and
 267 T-SAEs, may be complements rather than substitutes: TXCs perform well on synthetic temporal
 268 recovery and backtracking, while T-SAEs perform better on emergent misalignment and HH-RLHF
 269 decomposition. A more systematic temporal benchmark, across both synthetic processes and real-
 270 world tasks, would help clarify the capability profile of each architecture and is an important direction
 271 for further work.

272 Second, we focus on mapping the capability profile of temporal crosscoders and not a theoretical
273 account of when and why they outperform. This leads to three limitations. First, we cannot yet predict
274 in advance which real-world tasks will favor TXCs, T-SAEs, standard SAEs, or crosscoders. Second,
275 we do not derive optimal hyperparameters for TXCs, and we believe it likely that the flexibility
276 of the TXC framework will allow significant improvement on the results reported here. Third, an
277 intriguing question left unanswered by our work is whether the TSAE architectures can be effectively
278 incorporated into the TXC framework. This is particularly important since these two architectures
279 had complementary capability profiles on TempBench and combining the strengths of both into a
280 single architecture is an exciting possibility for future work.

281 Taken together, these gaps mean that significant work remains both to fully understand the TXC
282 framework and to fully exploit it. Mapping the regimes in which a temporal inductive bias pays off,
283 and combining the architectural strengths of TXCs and T-SAEs into a single dictionary, are in our
284 view the most exciting directions this paper points toward.

285 7 Discussion

286 Temporal crosscoders are currently best understood as a flexible addition to the sparse-dictionary
287 toolkit. Across TempBench, TXCs help when the target feature is supported by evidence spread over
288 nearby tokens: they recover hidden temporal variables in the synthetic benches and perform strongly
289 on sparse probing and backtracking detection/inducement. The negative cases are equally informative:
290 on emergent misalignment and HH-RLHF preference decomposition, temporal aggregation can
291 diffuse representations or amplify length-correlated artifacts. Thus, the right unit of explanation is
292 task-dependent: some features are token-local, while others are temporally extended.

293 Several directions follow naturally. First, TXCs currently use short fixed windows. Scaling to longer
294 contexts is a confirmed priority, and will likely require additional structure such as hierarchical
295 windows, low-rank methods, state-space compression, or tensor-network factorizations [Orús, 2014,
296 2019]. Second, transformer residual streams contain skip connections across layers, so temporal
297 feature discovery should eventually account for both sequence structure and residual-stream com-
298 position [Elhage et al., 2021]. Combining TXCs with cross-layer crosscoders is one direct way to
299 test this. Third, the architecture leaves substantial room for improvement: adaptive window lengths,
300 stronger contrastive losses, predictive objectives, and multi-scale training could separate genuine
301 temporal features from confounds such as response length and feature density. The central question
302 is therefore not whether temporal dictionaries are always better, but when model features are best
303 represented as local directions, temporal directions, or interactions between the two.

304 References

- 305 Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-
306 coders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*,
307 2023.
- 308 Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick
309 Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec,
310 Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina
311 Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christo-
312 pher Olah. Towards monosemanticity: Decomposing language models with dictionary learning.
313 *Transformer Circuits Thread*, 2023a. <https://transformer-circuits.pub/2023/monosemantic-features>.
- 314 Ekdeep Singh Lubana, Can Rager, Sai Sumedh R. Hindupur, Valérie Costa, Greta Tuckute, Oam
315 Patel, Sonia Krishna Murthy, Thomas Fel, Daniel Wurgaft, Eric J. Bigelow, Johnny Lin, Demba Ba,
316 Martin Wattenberg, Fernanda Viegas, Melanie Weber, and Aaron Mueller. Priors in time: Missing
317 inductive biases for language model interpretability. *arXiv preprint arXiv:2511.01836*, 2025. URL
318 <https://arxiv.org/abs/2511.01836>.
- 319 N. Chomsky. Three models for the description of language. *IRE Transactions on Information Theory*,
320 2(3):113–124, 1956. doi: 10.1109/TIT.1956.1056813.
- 321 Roger Levy. Expectation-based syntactic comprehension. *Cognition*, 106(3):1126–1177, 2008.

- 322 Yulia Lerner, Christopher J Honey, Lauren J Silbert, and Uri Hasson. Topographic mapping of a
323 hierarchy of temporal receptive windows using a narrated story. *Journal of Neuroscience*, 31(8):
324 2906–2915, 2011. doi: 10.1523/JNEUROSCI.3684-10.2011.
- 325 Usha Bhalla, Alex Oesterling, Claudio Mayrink Verdun, Himabindu Lakkaraju, and Flavio P. Calmon.
326 Temporal sparse autoencoders: Leveraging the sequential nature of language for interpretability,
327 2026. URL <https://arxiv.org/abs/2511.05541>.
- 328 Francesco Cagnetta and Matthieu Wyart. Towards a theory of how the structure of lan-
329 guage is acquired by deep neural networks. In A. Globerson, L. Mackey, D. Belgrave,
330 A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information*
331 *Processing Systems*, volume 37, pages 83119–83163. Curran Associates, Inc., 2024. doi:
332 10.52202/079017-2645. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2024/file/9740da1c07c7b451af14e11523f95271-Paper-Conference.pdf)
333 [2024/file/9740da1c07c7b451af14e11523f95271-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/9740da1c07c7b451af14e11523f95271-Paper-Conference.pdf).
- 334 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen,
335 Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L
336 Turner, Callum McDougall, Monte MacDiarmid, Alex Tamkin, Esin Durmus, Tristan Hume,
337 Francesco Mosconi, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson,
338 Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting
339 interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024. [https://transformer-](https://transformer-circuits.pub/2024/scaling-monosemanticity/)
340 [circuits.pub/2024/scaling-monosemanticity/](https://transformer-circuits.pub/2024/scaling-monosemanticity/).
- 341 Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya
342 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint*
343 *arXiv:2406.04093*, 2024. URL <https://arxiv.org/abs/2406.04093>.
- 344 Bart Bussmann, Patrick Leask, and Neel Nanda. BatchTopK sparse autoencoders, 2024. URL
345 <https://arxiv.org/abs/2412.06410>.
- 346 Tom Lieberum, Senthooan Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant
347 Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse
348 autoencoders everywhere all at once on gemma 2. In *Proceedings of the 7th BlackboxNLP*
349 *Workshop*, pages 278–300, 2024. URL <https://arxiv.org/abs/2408.05147>.
- 350 Subhash Kantamneni, Joshua Engels, Senthooan Rajamanoharan, Max Tegmark, and Neel Nanda.
351 Are sparse autoencoders useful? a case study in sparse probing, 2025. URL [https://arxiv.](https://arxiv.org/abs/2502.16681)
352 [org/abs/2502.16681](https://arxiv.org/abs/2502.16681).
- 353 David Chanin. Can saes recover true features from token embeddings? [https://colab.research.](https://colab.research.google.com/drive/1lLTvn6ksjjVZ0JtF4va_mkDGFf8TsZ7B?usp=sharing)
354 [google.com/drive/1lLTvn6ksjjVZ0JtF4va_mkDGFf8TsZ7B?usp=sharing](https://colab.research.google.com/drive/1lLTvn6ksjjVZ0JtF4va_mkDGFf8TsZ7B?usp=sharing), 2025. Colab
355 notebook.
- 356 Aleksandar Makelov, George Lange, and Neel Nanda. Towards principled evaluations of sparse au-
357 toencoders for interpretability and control, 2024. URL <https://arxiv.org/abs/2405.08366>.
- 358 Adam Karvonen, Benjamin Wright, Can Rager, Rico Angell, Jannik Brinkmann, Logan Smith,
359 Claudio Mayrink Verdun, David Bau, and Samuel Marks. Measuring progress in dictionary learning
360 for language model interpretability with board game models. *Advances in Neural Information*
361 *Processing Systems*, 37:83091–83118, 2024.
- 362 Jack Lindsey, Adly Templeton, Jonathan Marcus, Tom Conerly, Joshua Baston, and Chris Olah.
363 Sparse crosscoders for cross-layer features and model diffing. *Transformer Circuits Thread*, 2024.
364 URL <https://transformer-circuits.pub/2024/crosscoders/index.html>.
- 365 Julian Minder, Clément Dumas, Caden Juang, Bilal Chughtai, and Neel Nanda. Overcoming sparsity
366 artifacts in crosscoders to interpret chat-tuning, 2026. URL [https://arxiv.org/abs/2504.](https://arxiv.org/abs/2504.02922)
367 [02922](https://arxiv.org/abs/2504.02922).
- 368 Usha Bhalla, Alex Oesterling, Claudio Mayrink Verdun, Himabindu Lakkaraju, and Flavio P. Calmon.
369 Temporal sparse autoencoders: Leveraging the sequential nature of language for interpretability.
370 *arXiv preprint arXiv:2511.05541*, 2025. URL <https://arxiv.org/abs/2511.05541>.

- 371 Adam Karvonen, Can Rager, Johnny Lin, Curt Tigges, Joseph Bloom, David Chanin, Yeu-Tong Lau,
372 Eoin Farrell, Callum McDougall, Kola Ayonrinde, Demian Till, Matthew Wearden, Arthur Conmy,
373 Samuel Marks, and Neel Nanda. Saebench: A comprehensive benchmark for sparse autoencoders
374 in language model interpretability, 2025. URL <https://arxiv.org/abs/2503.09532>.
- 375 James P. Crutchfield and Karl Young. Inferring statistical complexity. *Physical Review Letters*, 63(2):
376 105–108, 1989.
- 377 James P. Crutchfield. Between order and chaos. *Nature Physics*, 8(1):17–24, 2012.
- 378 Adam S. Shai, Lucas Teixeira, Alexander Gilbert Oldenziel, Sarah Marzen, and Paul M. Riechers.
379 Transformers represent belief state geometry in their residual stream. In *Advances in Neural*
380 *Information Processing Systems*, 2024.
- 381 Adam S. Shai, Loren Amdahl-Culleton, Casper L. Christensen, Henry R. Bigelow, Fernando E. Rosas,
382 Alexander B. Boyd, Eric A. Alt, Kyle J. Ray, and Paul M. Riechers. Transformers learn factored
383 representations, 2026. URL <https://arxiv.org/abs/2602.02385>.
- 384 Francis Rhys Ward et al. Reasoning-finetuning repurposes latent representations in base models. In
385 *ICML 2025 Workshop on Actionable Interpretability*, 2025. arXiv:2507.12638.
- 386 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain,
387 Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with
388 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 389 Román Orús. A practical introduction to tensor networks: Matrix product states and projected
390 entangled pair states. *Annals of Physics*, 349:117–158, 2014. doi: 10.1016/j.aop.2014.06.013.
391 URL <https://doi.org/10.1016/j.aop.2014.06.013>.
- 392 Román Orús. Tensor networks for complex quantum systems. *Nature Reviews Physics*, 1:
393 538–550, 2019. doi: 10.1038/s42254-019-0086-7. URL <https://doi.org/10.1038/s42254-019-0086-7>.
- 395 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda
396 Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac
397 Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse,
398 Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A
399 mathematical framework for transformer circuits. Transformer Circuits Thread, 2021. URL <https://transformer-circuits.pub/2021/framework/index.html>. Published December 2021.
- 401 Trenton Bricken et al. Towards monosemanticity: Decomposing language models with dictionary
402 learning. <https://transformer-circuits.pub/2023/monosemantic-features>, 2023b.
403 Transformer Circuits Thread.
- 404 Constantin Venhoeff, Iván Arcuschin, Philip Torr, Arthur Conmy, and Neel Nanda. Understanding
405 reasoning in thinking language models via steering vectors. In *Reasoning and Planning for LLMs*
406 *@ ICLR2025*, 2025. URL <https://arxiv.org/abs/2506.18167>.
- 407 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
408 and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In
409 *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*,
410 2021. URL <https://arxiv.org/abs/2103.03874>.
- 411 Anthropic. Introducing Claude Haiku 4.5, Oct 2025. URL <https://www.anthropic.com/news/claude-haiku-4-5>. Accessed: 2026-05-06.
412

413 **A Model Parameters for each architecture**

414 This appendix specifies architecture and training choices held fixed across components. The only
 415 per-component free knobs are dictionary width d_{SAE} and sparsity budget k_{pos} , the expected num-
 416 ber of active latents per sequence position. Toy data uses $d_{\text{SAE}} = 40$; real-model widths are
 417 set by subject-model dimension. All locked choices are in `configs/locked_archs.yaml` and
 418 `configs/datasources.yaml`. Component-specific datasource, width, step count, k_{pos} grid, and
 419 seed tables appear in the component appendices.

420 **Common training pipeline.** All cells use the same trainer: Adam with $\text{lr} = 3 \times 10^{-4}$, batch
 421 size $B = 1024$, 1,000 warmup steps, and mixed precision, using bf16 on H100/H200 and
 422 fp16 with grad-scaler on A40. Toy-data cells run $n_{\text{steps}} \in \{8,000, 30,000\}$; real-LM cells run
 423 $n_{\text{steps}} \in \{10,000, 20,000\}$. Runs with $n_{\text{steps}} < 1000$ are discarded as smoke runs. Across the seven
 424 components, the paper uses $\approx 5,700$ trained checkpoints and $\approx 6,200$ leaderboard eval rows. Per-
 425 component counts and the full (arch, datasource, seeds) breakdown are given below. Subject-model
 426 activations are cached once per datasource and reused by all architectures at matched hookpoint. Spar-
 427 sity is matched by k_{pos} : per-token dictionaries use k_{pos} active latents per token, while a windowed
 428 architecture of length T uses

$$k_{\text{win}} = T k_{\text{pos}}.$$

429 For TXC-pro, matryoshka groups use

$$k_{\text{train}} = 5k_{\text{pos}}, \quad k_{\text{inference}} = 10k_{\text{pos}}.$$

430 **Anti-dead stack.** TopK SAE, T-SAE, TXC-base, and TXC-pro share the same anti-dead stack:
 431 (i) per-feature `num_tokens_since_fired` counters with a 10M-token dead threshold; (ii) AuxK
 432 loss with $\alpha_{\text{auxk}} = 1/32$, where the top-aux_k dead features re-reconstruct the residual; (iii) decoder
 433 unit-norm constraint per latent over (T, d_{in}) ; (iv) decoder-parallel gradient removal on W_{dec} ; and (v)
 434 geometric-median b_{dec} initialization on the first batch. This ports the Bhalla et al. [2025] anti-dead
 435 recipe into the TXC family. Bricken-style hard resampling [Bricken et al., 2023b] is off by default
 436 and used only in C6 to match the Qwen-7B medical setup; see section D.

437 **TXC-base.** TXC-base is the acausal temporal crosscoder from section 3 with TopK sparsity and
 438 the full anti-dead stack. The locked window length is $T = 5$. For each $\tau \in [0, T)$, encoder slabs

$$W_{\text{enc}}^{(\tau)} \in \mathbb{R}^{d_{\text{in}} \times d_{\text{SAE}}}$$

439 sum into one shared latent vector $\mathbf{u}_t \in \mathbb{R}^{d_{\text{SAE}}}$ per window. BatchTopK is applied over the latent axis
 440 with $k_{\text{win}} = T k_{\text{pos}}$. Decoder slabs

$$W_{\text{dec}}^{(\tau)} \in \mathbb{R}^{d_{\text{SAE}} \times d_{\text{in}}}, \quad b_{\text{dec}}^{(\tau)} \in \mathbb{R}^{d_{\text{in}}},$$

441 reconstruct the full T -window. TXC-base uses no matryoshka objective, no contrastive loss, and no
 442 Bricken resampling. The free knobs are k_{pos} and d_{SAE} . C3 also trains TXC-base at $T \in \{10, 20\}$ to
 443 probe temporal-context length; see section C.1.

444 **TXC-pro.** TXC-pro adds three choices to TXC-base. First, the subseq encoder trains with $T_{\text{max}} =$
 445 10 position slabs while sampling $t_{\text{sample}} = 5$ positions per step; all 10 positions are used at evaluation.
 446 Second, Matryoshka training uses $H = 8$ nested feature groups, each reconstructing from the first

$$G d_{\text{SAE}} / 8, \quad G = 1, \dots, 8$$

447 features. Third, a multi-distance contrastive loss applies inverse-distance-weighted InfoNCE at shifts
 448 $\Delta \in \{1, 2\}$. The matryoshka head is disabled at toy scale in C1 and C2 by setting $h_{\text{size}} = d_{\text{SAE}} = 40$,
 449 since 5 prefix levels of 8 features do not form a useful hierarchy and the smallest prefix would fall
 450 below $k_{\text{train}} = 5k_{\text{pos}}$ at high k_{pos} . The free knobs are again k_{pos} and d_{SAE} .

451 **Baselines.** **TopK SAE** [Gao et al., 2024] is a per-token dictionary with TopK at k_{pos} and the
 452 anti-dead stack. **Stacked SAE**, used in C1, C2, and C7, is a bank of T independent TopK SAEs, one
 453 per position in a T -window, with total window budget $k_{\text{win}} = T k_{\text{pos}}$. It isolates temporal aggregation
 454 from cross-position weight sharing. **T-SAE** [Bhalla et al., 2025] is a faithful port of the referenced
 455 work’s implementation: Matryoshka BatchTopK with $h_{\text{frac}} = 0.20$ high-level features, AuxK loss,

456 geometric-median b_{dec} initialization, decoder unit-norm projection, raw-dot temporal contrastiveA
 457 loss at $\alpha = 1.0$, and threshold-based inference. Its locked $d_{\text{SAE}} = 16,384$ matches the original
 458 Gemma-2-2B setting. **TFA** [Lubana et al., 2025] is the temporal feature analyzer with 4 attention
 459 heads; the no-position-embedding variant is the default, and TFA-pos is reported on synthetic data.
 460 **MLC** [Lindsey et al., 2024] is a 5-layer multi-layer crosscoder over the residual-stream layer axis,
 461 centered at the component hookpoint. It uses layers 11–15 for Gemma-2-2B-IT C3 and layers 8–12
 462 for Llama-3.1-8B C7. **SAE-arditi** is used only as the C6 baseline to match the Qwen-medical setup,
 463 with $d_{\text{SAE}} = 32,768$ and $k_{\text{pos}} = 128$; see section D.

464 **Locked versus free choices.** The locked TXC-base and TXC-pro definitions above are the only
 465 architectural commitments. Across all seven components, only k_{pos} and d_{SAE} vary by component. We
 466 do not hill-climb architecture hyperparameters such as T , matryoshka group count, contrastive shifts,
 467 AuxK weight, or dead threshold. Training-time augmentations are allowed only when documented in
 468 the component appendix; the only such opt-in is Bricken resampling for C6.

469 B Synthetic Benchmarks: Ground-Truth Correspondences and Probe Setup

470 The denoising metrics in section 4 (single-latent correlation and linear-probe R^2) rely on two
 471 distinct correspondences between ground-truth objects and the architecture’s dictionary: one fixed by
 472 construction in the data generator, one inferred post-hoc by decoder geometry. This appendix makes
 473 both explicit so that the local-vs-global scatters in section 4 can be interpreted unambiguously.

474 B.1 Construction-time correspondence (feature index \leftrightarrow chain)

475 The Denoising bench generator (`temp_bench.data.toy.markov:markov_chain_support`)
 476 draws N orthogonal feature directions $\mathbf{f}_i \in \mathbb{R}^d$ as the rows of a QR-orthogonalised Gaussian
 477 matrix, then runs N independent two-state Markov chains $s_{i,t} \in \{0, 1\}$ in parallel and emits noisy
 478 firings $a_{i,t} \mid s_{i,t} \sim \text{Bernoulli}(p_A [s_{i,t}=0] + p_B [s_{i,t}=1])$. The feature index $i \in \{1, \dots, N\}$ simulta-
 479 neously labels three objects: the orthogonal direction \mathbf{f}_i , the hidden state $s_{i,t}$, and the noisy emission
 480 $a_{i,t}$. The activation

$$\mathbf{x}_t = \sum_{i=1}^N a_{i,t} m_{i,t} \mathbf{f}_i$$

481 threads this index through the residual-stream construction, so the feature \leftrightarrow emission \leftrightarrow hidden-state
 482 mapping is determined entirely by the generator. This is the ground-truth correspondence the metrics
 483 evaluate against.

484 The Coupling bench (`temp_bench.data.toy.coupled_noisy:coupled_noisy_hmm`) keeps the
 485 same construction-time logic for the $M = 20$ emissions $a_{j,t}$ and adds a separate index $k \in$
 486 $\{1, \dots, K\}$ with $K = 10$ for the hidden chain directions $\mathbf{f}_k^{\text{hid}}$. A binary coupling matrix
 487 $C \in \{0, 1\}^{M \times K}$ specifies which subset of hidden chains drives each emission via OR-gate fir-
 488 ing, so the construction-time correspondence has two layers (chain index k , emission index j) instead
 489 of one.

490 B.2 Post-hoc correspondence (feature direction \leftrightarrow dictionary atom)

491 For an autoencoder with decoder columns $\mathbf{w}_j \in \mathbb{R}^d$, $j \in \{1, \dots, d_{\text{sae}}\}$, the matched latent for
 492 ground-truth feature i is the column with the largest absolute cosine,

$$j^*(i) = \arg \max_j \left| \frac{\langle \mathbf{w}_j, \mathbf{f}_i \rangle}{\|\mathbf{w}_j\| \|\mathbf{f}_i\|} \right|.$$

493 The matching is computed independently for each feature, so two ground-truth features may in
 494 principle pick the same latent. We do not enforce uniqueness; this matches the wasteland-source
 495 convention in `denoising_probes.py`.

496 B.3 Probe definitions

497 Given the eval set of activations \mathbf{x}_t , every architecture produces a per-token latent code $\mathbf{u}(t) \in \mathbb{R}^{d_{\text{sae}}}$.
 498 For per-token architectures (TopK SAE, T-SAE, TFA-pos) the encoder runs directly on each token.

499 For window-level architectures (TXC, Stacked) the encoder produces one code per T -window; we
 500 broadcast that code to all T positions in the window and average over the T overlapping windows
 501 that touch each position, so every (seq, t) pair ends up with one d_{sae} -dim vector.

502 **Single-latent correlation.** For each ground-truth feature i , take the matched latent $j^*(i)$ defined
 503 above and compute Pearson correlations over all (seq, t) pairs,

$$\begin{aligned}\bar{r}_{\text{local}}(i) &= \text{Corr}(u_{j^*(i)}(t), a_{i,t}), \\ \bar{r}_{\text{global}}(i) &= \text{Corr}(u_{j^*(i)}(t), s_{i,t}),\end{aligned}$$

504 then average across the N features. This is a *strict* test: the dictionary direction must align with \mathbf{f}_i
 505 and the code dimension along that direction must track the hidden state.

506 **Linear probe.** For each feature i fit two ridge regressions ($\alpha = 1$, 80/20 train/test split) using the
 507 full code $\mathbf{u}(t)$,

$$\begin{aligned}\hat{a}_{i,t} &= \mathbf{w}_{\text{loc}}^{(i)\top} \mathbf{u}(t) + b_{\text{loc}}^{(i)}, \\ \hat{s}_{i,t} &= \mathbf{w}_{\text{glob}}^{(i)\top} \mathbf{u}(t) + b_{\text{glob}}^{(i)},\end{aligned}$$

508 and report $R_{\text{local}}^2(i)$ and $R_{\text{global}}^2(i)$ on the held-out 20%, averaged across the N features. This is a
 509 *looser* test: information may be diffused across many latents and the probe still recovers it.

510 The local-vs-global scatters in section 4 use both metrics on the same population of cells. Points
 511 above the $y = x$ diagonal denote: the latent code tracks the hidden state better than the noisy
 512 emission.

513 C Sparse Probing: Full Details

514 Sparse probing measures whether a dictionary exposes features useful for simple classifiers of
 515 human-interpretable text concepts. Standard probing panels include both token-local properties, such
 516 as punctuation and part of speech, and broader passage-level properties, such as topic, sentiment,
 517 and named entities. This makes sparse probing a useful static test of whether temporal-window
 518 dictionaries expose concepts that per-token dictionaries fragment.

519 C.1 Architectures and training configuration

520 Six base architectures, each trained on the same Gemma-2-2B-IT layer-13 residual-stream activation
 521 cache (FineWeb 24K sequences \times 128 tokens, fp16): a per-token TopK SAE Gao et al. [2024], the
 522 temporal SAE of Bhalla et al. [2025], the temporal feature analyzer (TFA) of Lubana et al. [2025], a
 523 multi-layer crosscoder (MLC) on layers 11–15 Lindsey et al. [2024], TXC-base ($T = 5$, $k_{\text{pos}} = 20$),
 524 and TXC-pro ($T_{\text{max}} = 10$, $t_{\text{sample}} = 5$, $k_{\text{pos}} = 20$, eight Matryoshka groups, contrastive shifts
 525 $\Delta \in \{1, 2\}$). All cells train for $n_{\text{steps}} = 20,000$ at $d_{\text{SAE}} = 18,432$. The matched-sparsity invariant
 526 fixes the expected per-token L_0 at 20 across architectures; each arch’s internal TopK budget scales
 527 with the unit it operates on (TopK SAE: $k_{\text{win}} = 20$; T-SAE BatchTopK: $20 \cdot B \cdot T$; TXC-base (T):
 528 $k_{\text{win}} = 20 \cdot T$; TXC-pro: $k_{\text{train}} = 100$, $k_{\text{inference}} = 200$; MLC ($L = 5$): $k_{\text{win}} = 100$; TFA: $20 \cdot T_{\text{tfa}}$).
 529 Within-component fairness is enforced by this matched-sparsity convention rather than by uniform
 530 per-step token throughput. We additionally train TXC-base at $T \in \{10, 20\}$ as a within-arch probe of
 531 temporal-context length.

532 C.2 Probing protocol

533 For every probe task we cache the last $S = 32$ token activations of the subject model on each prompt.
 534 Tokenisation is right-padded then per-example resliced into a 32-frame so the real tokens occupy
 535 positions $[\text{first_real}_i, S - 1]$ and zero padding sits at the start; the per-example `first_real` index
 536 is stored so downstream pooling can mask the padding region. Per-token architectures encode the
 537 full (N, S, d_{in}) tensor and mean-pool over the token axis with positions $k < \text{first_real}_i$ masked.
 538 Window architectures slide a width- T window across the S -frame at stride 1 ($n_{\text{windows}} = S - T + 1$
 539 per prompt; with $S = 32$: $T = 5 \rightarrow 28$, $T = 10 \rightarrow 23$, $T = 20 \rightarrow 13$), encode each window
 540 to one (d_{SAE}) latent, then mean-pool over the window axis with windows whose left edge lies in

padding masked. The output is a single (d_{SAE}) feature vector per prompt for every architecture. Top- S features are selected on the train split only by $|\mu_{y=1} - \mu_{y=0}|$ (no CV / no peeking at the test set during selection), and the probe is scikit-learn LogisticRegression(penalty='l1', solver='liblinear', C=1.0, max_iter=1000, random_state=0) fit on the selected k_{feats} features. The per-task readout is ROC-AUC on the task’s held-out test split; the headline is the unweighted mean across the SAEBench panel, averaged over 3 training seeds.

547 C.3 SAEBench task suite

548 The probing panel is the upstream SAEBench composition fixed by
 549 chosen_classes_per_dataset in sae_bench/sae_bench_utils/dataset_info.py, to-
 550 tallying 36 binary one-vs-rest tasks across 8 datasets:

dataset	classes	n
LabHC/bias_in_bios_class_set1	professions 0, 1, 2, 6, 9	5
LabHC/bias_in_bios_class_set2	professions 11, 13, 14, 18, 19	5
LabHC/bias_in_bios_class_set3	professions 20, 21, 22, 25, 26	5
canrager/amazon_reviews_mcauley_1and5	categories 1, 2, 3, 5, 6	5
canrager/amazon_reviews_mcauley_1and5_sentiment	1.0-vs-rest, 5.0-vs-rest	2
codeparrot/github-code	C, Python, HTML, Java, PHP	5
fancyzhx/ag_news	world, sports, business, scitech	4
Helsinki-NLP/europarl	en, fr, de, es, nl	5
SAEBench total		36

551 Two cross-token tasks (winogrande_correct_completion and wsc_coreference) round out a
 552 38-task panel; the main-text AUC summary averages over all 38 tasks. Per-task AUCs are persisted
 553 in the leaderboard for downstream diagnostic queries.

554 C.4 Sweep aggregation: $\overline{\text{AUC}}$

555 The per-arch summary metric reported in fig. 3a compresses the k_{feats} sweep into a single number.
 556 Let $\text{AUC}_a(k)$ be the mean SAEBench AUC of arch a at $k_{\text{feats}} = k$ averaged over seeds. We define

$$\overline{\text{AUC}}(a) = \frac{1}{\log_2(k_{\text{max}}/k_{\text{min}})} \sum_{i=0}^{n-2} \frac{\text{AUC}_a(k_i) + \text{AUC}_a(k_{i+1})}{2} (\log_2 k_{i+1} - \log_2 k_i), \quad (1)$$

557 where $\{k_i\}_{i=0}^{n-1} = \{5, 10, 20, 40, 80, 160, 320, 640\}$. The grid is geometric so $\log_2 k_{i+1} - \log_2 k_i = 1$
 558 for every i , and eq. (1) reduces to the unweighted trapezoidal mean over $\log_2 k_{\text{feats}}$. Seed- σ on the
 559 summary metric is computed by propagating per-cell seed- σ through the trapezoidal weights.

560 C.5 Per-task breakdown at $k_{\text{feats}} = 160$

561 Figure 6 shows the architecture \times task ROC-AUC matrix at the peak-AUC sparsity $k_{\text{feats}} = 160$,
 562 averaged over the 3 training seeds. Most cells in the temporal family saturate above 0.90 on the
 563 in-distribution datasets (Bias in Bios, Amazon, GitHub, AG News, Europarl), with TopK SAE
 564 consistently ~ 0.01 – 0.04 behind. Two patterns are visible: TFA underperforms by ~ 0.05 – 0.15 AUC
 565 on every in-distribution group, confirming that the gap reported in fig. 3a is broad rather than driven
 566 by a small subset of tasks; and the two cross-token tasks (winogrande, wsc) sit at 0.13–0.40 for
 567 every architecture, reflecting that a per-token mean-pool encoder cannot represent the long-range
 568 coreference contrast these tasks require regardless of the SAE inductive bias (section C.7).

569 C.6 Full AUC curve with TFA

570 Figure 7 shows the same per-architecture probe sweep as the right panel of fig. 3 but with TFA
 571 included on the same axes. We separate TFA into the appendix because its AUC range (~ 0.65 at
 572 $k_{\text{feats}} = 5$, climbing to ~ 0.85 at $k_{\text{feats}} = 640$) compresses the y -axis on the main-text panel and
 573 obscures the ~ 0.01 – 0.02 AUC ordering among the temporal family and the per-token TopK SAE
 574 baseline. TFA does not finish saturating within the $k_{\text{feats}} \leq 640$ sweep; the gap between TFA and the

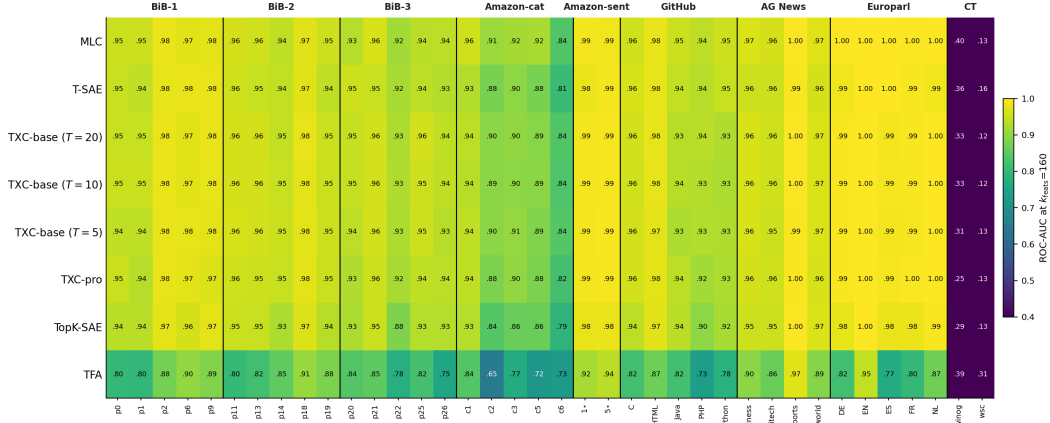


Figure 6: Per-architecture, per-task ROC-AUC at $k_{\text{feats}} = 160$ on the 38-task SAEBench panel, averaged over 3 training seeds. Rows ordered top-down by overall ranking; columns grouped by source dataset (BiB-1/2/3 = bias_in_bios sets, Amazon-cat = 5-class category, Amazon-sent = 1.0-vs-rest / 5.0-vs-rest, CT = the two cross-token coreference tasks). Cell labels are the AUC values with the leading “0.” dropped.

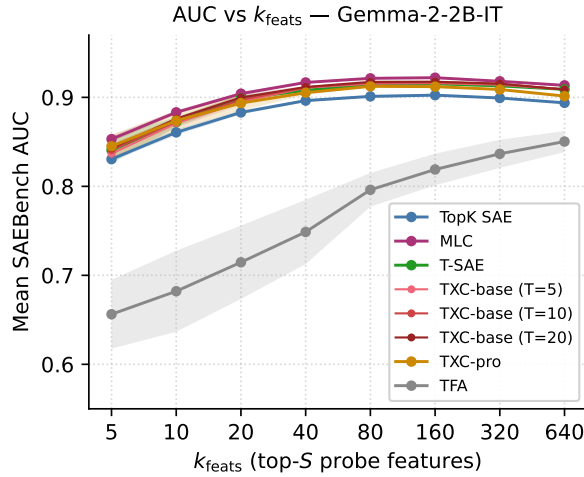


Figure 7: Full mean SAEBench AUC vs k_{feats} on Gemma-2-2B-IT for every architecture, including TFA. Lines are mean over 3 seeds; shaded band is $\pm 1\sigma$ across seeds. The temporal family + TopK SAE cluster sits at the top (saturating near $k_{\text{feats}} = 160$ at $\text{AUC} \approx 0.90\text{--}0.92$); TFA underperforms by 0.05–0.18 AUC across the sweep and is still climbing at the largest k_{feats} tested.

575 temporal family at the largest k_{feats} is well outside the per-cell seed- σ (≤ 0.02 for TFA, ≤ 0.002 for
 576 every other architecture).

577 C.7 Caveats

578 • **Cross-token tasks.** winogrande_correct_completion and wsc_coreference require
 579 resolving a referent that lives many tokens from the variant word; a per-token mean-pool
 580 encoder cannot represent the relevant contrast, and every locked-arch cell scored 0.40–0.50
 581 AUC on them irrespective of architecture (fig. 6). Including them shifts every $\overline{\text{AUC}}$ down
 582 by ~ 0.027 uniformly with no rank flips except TFA, which is already low.

583 • **TFA at paper-faithful $B = 32$.** TFA’s $\overline{\text{AUC}} = 0.764$ underperforms the locked set by
 584 0.07–0.18 AUC across the sweep; we attribute this to the small batch size required by the

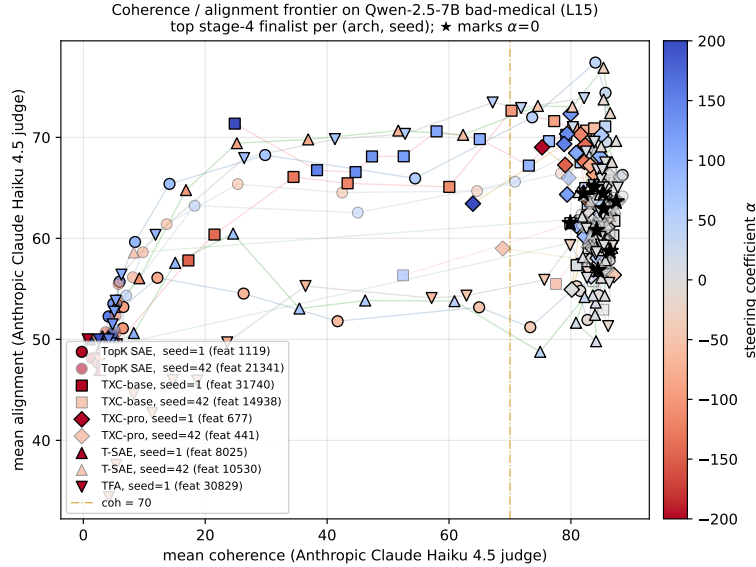


Figure 8: Coherence/alignment Pareto frontier for the Qwen-7B medical-misalignment cells. Each curve follows one architecture/seed across steering magnitudes; the vertical extent to the right of the $\text{coh} = 70$ threshold defines the dynamic range used in fig. 5a.

585 (B, T, d_{SAE}) -sized attention tensor at $d_{\text{SAE}} = 18,432$. A fairer comparison at $B = 1024$
 586 would require re-architecting the attention head and is left for follow-up.

587 • **Single subject model.** The headline is on Gemma-2-2B-IT only; replication on the BASE
 588 side and on a different subject model is left for follow-up.

589 D Emergent Misalignment Case Study: Full Steering Grid

590 The two summary subfigures in fig. 5 collapse, for each architecture, the full Wang stage-4 α
 591 frontier into a single coherence-floored alignment dynamic range (fig. 5a) and a single sparse-probe
 592 PR-AUC (fig. 5b). Figure 8 replots the same data as a coherence/alignment Pareto scatter (one
 593 curve per (architecture, seed), each curve following the top stage-4 finalist’s full α frontier), making
 594 the dynamic range visible as the vertical extent of points to the right of the $\text{coh} = 70$ threshold.
 595 Figure 9 shows the same data in unaggregated form: each panel is one (architecture, seed) cell on
 596 Qwen-2.5-7B-Instruct + bad-medical LoRA, with the mean alignment of all three Wang stage-4
 597 finalists plotted against α , the top-finalist mean coherence overlaid as a dashed line, and the $\alpha = 0$
 598 unsteered baseline marked in green. The shaded green band marks the $\text{coh} \geq 70$ region; red and
 599 blue triangles mark the extrema used in eq. (2). The dense per-finalist α grid combines the 27-point
 600 Wang stage-4 grid (densely sampled in $|\alpha| \in [1, 10]$ with $\alpha = \pm 100$ at the tails) with the extended
 601 $|\alpha| \in \{110, 120, 130, 140, 150, 200\}$ sweep run on the top finalist of each cell. The dense extension is
 602 what reveals coherence-preserving alignment peaks at $|\alpha| > 10$ that Han’s canonical grid missed (e.g.
 603 SAE-arditi seed 1, top finalist feat 1119, peaks at $\alpha = +40$ with align 77.4, $\text{coh} > 70$). All judge
 604 transcripts (per-rollout completions and Anthropic Claude Haiku 4.5 align/coh scores) are persisted in
 605 `judge_outputs.jsonl` (canonical) and `judge_outputs_extended.jsonl` (extension) for post-
 606 deadline κ validation.

607 **Alignment dynamic range.** For each architecture/seed cell c , let \mathcal{A}_c be the set of steering magni-
 608 tudes whose judged coherence is at least 70. We summarize steering control by the coherence-floored
 609 alignment range

$$\Delta_{\text{align}}(c) = \max_{\alpha \in \mathcal{A}_c} \text{align}(c, \alpha) - \min_{\alpha \in \mathcal{A}_c} \text{align}(c, \alpha), \quad \mathcal{A}_c = \{\alpha : \text{coh}(c, \alpha) \geq 70\}. \quad (2)$$

610 This is the statistic plotted in fig. 5a.

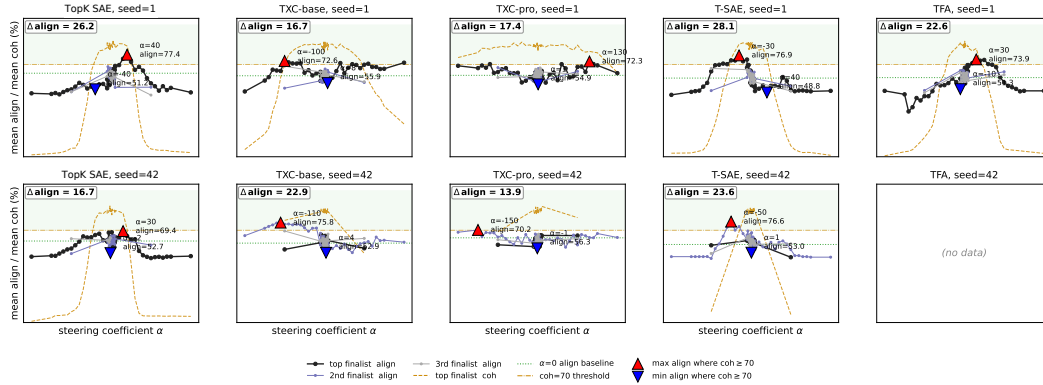


Figure 9: Full steering sweep for the Qwen-7B medical-misalignment cells. Panels show alignment across steering magnitudes, with coherence overlaid for the top finalist; extrema inside the $\text{coh} \geq 70$ region determine the summarized dynamic range.

611 D.1 Decoder-norm calibration ($1/\sqrt{T}$)

612 When a TXC reconstructs T positions from a single shared latent, each position-specific decoder
 613 row carries only a fraction of the per-position reconstructive mass that a per-token SAE decoder
 614 row carries; under the BatchTopK normalisation convention used in section 3 this leaves each TXC
 615 decoder row at ℓ_2 norm $\approx 1/\sqrt{T}$ relative to the unit-norm convention of a per-token SAE decoder
 616 row. A magnitude α applied to a TXC mined direction therefore produces an effective residual-stream
 617 perturbation that is a factor of $1/\sqrt{T}$ smaller than the same α applied to a per-token SAE direction.
 618 To make the cross-architecture α -grid comparable, we rescale the TXC’s grid by \sqrt{T} , so that the
 619 same α corresponds to the same ℓ_2 -norm intervention across architectures.

620 We confirm this empirically across the canonical c6 cells: the mean ℓ_2 -norm of mined TXC-base
 621 ($T = 5$) decoder rows is 0.443, against the predicted $1/\sqrt{5} = 0.4472$; for TXC-pro ($T = 10$) the
 622 mean is 0.310, against $1/\sqrt{10} = 0.3162$. Both match the prediction within $\sim 1.5\%$, justifying the
 623 \sqrt{T} rescaling without per-cell calibration.

624 E HH-RLHF Preference Decomposition: Full Details

625 Figure 10 collapses the per-feature decomposition for each architecture into a top-20 tier-count panel
 626 and a semantic-mass-share panel. This appendix supplies the methodology behind that summary
 627 and lists the top-5 autointerpreted features per architecture (table 1). The qualitative reading of
 628 the architecture comparison is concrete in that table: T-SAE at the paper-faithful $k = 20$ surfaces
 629 alignment-relevant labels (*clarification questions and requests for explanation, expressions of limita-*
 630 *tion or inability in responses, AI capability limitations or disclaimers*), while TXC at $T = 5$,
 631 $k_{\text{win}} = 500$ surfaces sentence-level discourse markers (*conversational filler or discourse markers,*
 632 *conversational responses and speech patterns*) whose length correlations are exactly the spurious
 633 channel the paper-faithful contrastive-loss training is tuned to suppress.

634 E.1 Response masking and ranking metric

635 For each (*chosen*, *rejected*) pair, we tokenise the two strings independently with `max_length=`
 636 `256` and identify the differing assistant response by a character-level longest common prefix between
 637 the two raw strings; we then map character offsets back to token indices via the tokenizer’s offset
 638 map to produce a per-side `response_mask` that selects only the differing-response tokens. Layer-12
 639 residual-stream activations are captured under this mask, and per-feature means are taken across
 640 all valid response-token positions in the 1000 pairs. Features are ranked by the signed metric
 641 $\mu_{\text{rejected}} - \mu_{\text{chosen}}$, matching Bhalla et al. [2025, §4.5]; negative values flag features that activate
 642 more strongly on *chosen* responses, positive values on *rejected*. Tier classification thresholds $|r| < 0.2$

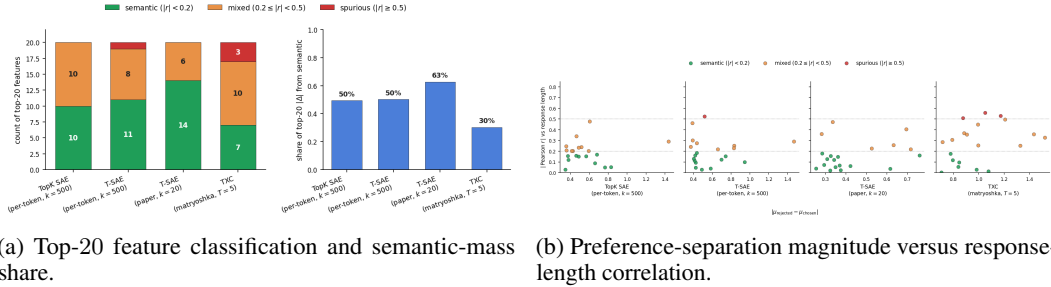


Figure 10: HH-RLHF preference decomposition for Gemma-2-2B layer 12. Features are ranked by $\mu_{\text{rejected}} - \mu_{\text{chosen}}$ and classified by correlation with response-length difference.

643 (semantic) and $|r| \geq 0.5$ (length-spurious) follow the paper. Our absolute response-length means
 644 (rejected 36.23, chosen 28.57 tokens) are lower than the paper’s reported 49.24 and 37.84 because
 645 of tokenizer and truncation differences across HH-RLHF dataset versions; the paired- t -test signal
 646 direction and significance ($p = 9.76 \times 10^{-10}$) match the paper’s $p \approx 9 \times 10^{-10}$ to leading digit.

647 For each feature j , define its preference separation

$$\Delta\mu_j = \mu_{\text{rejected},j} - \mu_{\text{chosen},j},$$

648 where means are taken over response-token positions only. We also compute the length correlation

$$r_j = \text{Corr}_{\text{pairs}}(\text{act}_{\text{rejected},j} - \text{act}_{\text{chosen},j}, \ell_{\text{rejected}} - \ell_{\text{chosen}}).$$

649 Features are classified as semantic if $|r_j| < 0.2$, mixed if $0.2 \leq |r_j| < 0.5$, and length-spurious if
 650 $|r_j| \geq 0.5$. For a top- K set \mathcal{T}_K , the semantic-mass share is

$$\text{SemMass}(\mathcal{T}_K) = \frac{\sum_{j \in \mathcal{T}_K: |r_j| < 0.2} |\Delta\mu_j|}{\sum_{j \in \mathcal{T}_K} |\Delta\mu_j|}.$$

651 E.2 Autointerpretation setup

652 For each top- K feature per architecture, we identify the 5 max-activating positions across the cohort
 653 (restricted to response tokens, so the label captures the chosen-vs-rejected differential signal rather
 654 than generic prompt activations), decode a symmetric token window around each firing position with
 655 the firing token marked by ****double asterisks****, and send the 5 contexts in a single Bills-et-al-
 656 style ? prompt to Claude Haiku 4.5 (claude-haiku-4-5-20251001, max output 80 tokens, default
 657 sampling). The prompt asks for a single short concept label (3–12 words), no preamble, no markdown.
 658 We use Haiku rather than Sonnet because the task is descriptive labeling rather than judgment; the
 659 lower cost lets us label the full top-50 per architecture for \sim \$0.25 total. All call records, contexts,
 660 and labels are persisted to the per-architecture `top_features.json` files referenced above.

661 F Backtracking Case Study: Full Details

662 This appendix gives implementation details for section 5.2.

663 F.1 Architectures and training configuration

664 We compare six architectures trained on the same Llama-3.1-8B residual-stream activation cache at
 665 layer 10, with $d_{\text{SAE}} = 32,768$, $n_{\text{steps}} = 300,000$, and batch size 1024: a per-token TopK SAE Gao
 666 et al. [2024], the temporal SAE of Bhalla et al. [2025], the temporal feature analyzer (TFA) of Lubana
 667 et al. [2025], a multi-layer crosscoder (MLC) Lindsey et al. [2024], TXC-base ($T = 5$, $k_{\text{pos}} = 20$),
 668 and TXC-pro ($T_{\text{max}} = 10$, $t_{\text{sample}} = 5$, $k_{\text{pos}} = 20$, eight Matryoshka groups, and contrastive shifts
 669 $\Delta \in \{1, 2\}$). We also train TXC-base and TXC-pro at batch size 256; these cells are reported in
 670 section F.13.

rank	r	$\mu_{\text{rejected}} - \mu_{\text{chosen}}$	autointerp label
<i>TopK SAE (per-token, k = 500)</i>			
0	+0.29	+1.439	objects or entities being directly affected by actions
1	-0.05	+0.830	physical objects or locations used in scenarios
2	-0.05	-0.763	information or data items that are absent or unavailable
3	-0.17	+0.680	illegal drugs and controlled substances
4	+0.09	+0.660	direct objects of transitive verbs or noun phrases
<i>T-SAE (per-token, k = 500)</i>			
0	+0.29	+1.468	pronouns and discourse connectives in explanatory contexts
1	-0.10	-0.944	negation or negative polarity markers
2	-0.25	-0.833	seeking confirmation or agreement on statements
3	-0.23	-0.828	hedging language and softening expressions
4	-0.15	-0.743	uncertainty markers in clarification requests
<i>T-SAE (paper-faithful, k = 20)</i>			
0	-0.16	-0.759	clarification questions and requests for explanation
1	-0.22	-0.716	clarification requests and confused responses
2	-0.41	-0.695	expressions of limitation or inability in responses
3	+0.26	+0.633	verb forms in instructional or procedural contexts
4	-0.04	-0.617	AI capability limitations or disclaimers
<i>TXC (matryoshka, T = 5, k_{win} = 500)</i>			
0	+0.33	+1.519	nouns denoting objects or entities being modified
1	-0.36	-1.433	discourse markers indicating uncertainty or topic shifts
2	-0.25	-1.327	clarification requests or seeking more information
3	-0.50	-1.207	conversational filler or discourse markers
4	-0.53	-1.172	conversational responses and speech patterns

Table 1: Top-5 features per architecture, ranked by $|\mu_{\text{rejected}} - \mu_{\text{chosen}}|$ on the 1000-pair harmless-base cohort. Sign of the diff indicates the side the feature activates more strongly on (negative \rightarrow chosen, positive \rightarrow rejected). Length-Pearson r classifies each feature as semantic ($|r| < 0.2$), mixed, or length-spurious ($|r| \geq 0.5$). Autointerp labels are produced by Claude Haiku 4.5 (section E.2).

671 F.2 Reproduction setup

672 We follow the dataset recipe of Venhoff et al. [2025], as reproduced in Ward et al. [2025, Ap-
673 pendix C]. The procedure has three steps: generate 300 math prompts with Claude Sonnet 3.7,
674 sample reasoning traces from DeepSeek-R1-Distill-Llama-8B, and label each sentence with an
675 LLM judge. The prompts cover ten categories, with 30 prompts per category: `basic_logic`,
676 `geometry`, `probability`, `arithmetic`, `counting`, `number_theory`, `set_theory`, `sequences`,
677 `inequalities`, and `algebra_word_problems`. We use 280 prompts for the `dom` split and hold out
678 20 for `eval`. Following our inducement-evaluation pipeline, we use Claude Sonnet 4.6 rather than
679 GPT-4o for sentence labels. Labels are binary, with field `is_backtracking`.

680 The steering vector is the difference-of-means direction for `is_backtracking`, computed from
681 base-model residual activations at layer 10. Activations are taken from token offsets -13 through
682 -8 relative to the first token of each labeled backtracking sentence. We use the same layer, offset
683 window, and base-model activation convention as Ward et al. [2025, Sec. 3.1, Appendix B.1].

684 F.3 Cohort construction

685 The evaluation cohort contains 61 MATH-500 questions Hendrycks et al. [2021]. It has two strata: 31
686 questions where the unsteered model gives a parsable but incorrect boxed answer, and 30 questions
687 where it gives a parsable correct boxed answer. We drop questions with no boxed answer within the
688 token budget, since both the inducement and answer-change metrics require a parsable continuation.
689 Questions are drawn deterministically from a locked `flip_matrix.parquet` file, so all architectures
690 are evaluated on the same questions in the same order. Question identifiers use the MATH-500
691 `unique_id` format, for example `test/algebra/1184.json`.

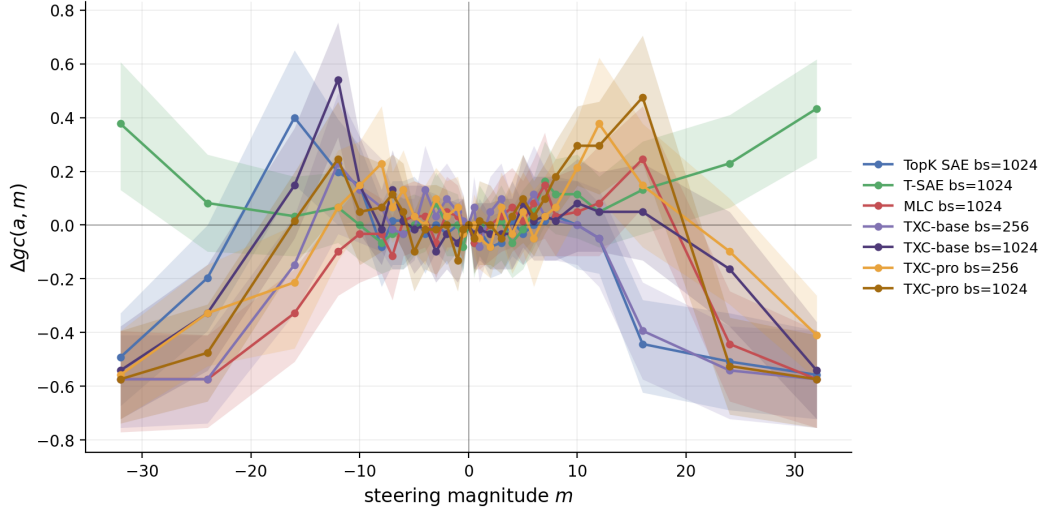


Figure 11: Per-architecture inducement curves $\Delta gc(a, m)$. Curves are per-question baseline-corrected at $m = 0$; shaded bands are bootstrap 95% CIs over the 61 cohort questions.

692 F.4 Magnitude grid

693 We evaluate the 25-point grid

$$\mathcal{M} = \{0, \pm 0.5, \pm 1, \pm 2, \pm 3, \pm 4, \pm 5, \pm 6, \pm 7, \pm 8, \pm 10, \pm 12, \pm 16\}.$$

694 The grid is denser between ± 0.5 and ± 8 , where the inducement curves change most quickly, and
 695 coarser at larger magnitudes. The $m = 0$ point is the no-intervention baseline used in section F.5.
 696 Figure 11 plots the per-architecture $\Delta gc(a, m)$ curves on this grid plus the extended follow-up at
 697 $m \in \{\pm 24, \pm 32\}$.

698 F.5 Inducement metric

699 For each architecture a , question q , and steering magnitude m , let $gc(a, q, m)$ be the Sonnet judge
 700 count of genuine backtracking events in the continuation. The main inducement statistic subtracts the
 701 shared unsteered cut-and-continue floor questionwise:

$$\Delta gc(a, m) = \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} (gc(a, q, m) - gc(a, q, 0)),$$

702 where \mathcal{Q} is the locked 61-question MATH-500 cohort. Positive Δgc means steering induces more
 703 genuine backtracking than the unsteered continuation protocol on the same questions.

704 F.6 cut25 continuation protocol

705 For each question q , let t_q^{full} be the token sequence of the unsteered reproduction. The prefix is the
 706 first $\lfloor 0.25 |t_q^{\text{full}}| \rfloor$ tokens. Conditioned on this prefix, the reasoning model samples a new continuation
 707 at temperature 1.0 with a 1024-token budget. During this continuation, the steering hook adds mv_{feat}
 708 to the layer-10 residual stream at each generated position. Here v_{feat} is the architecture-specific mined
 709 direction, rescaled to match the ℓ_2 norm of the dom-base-union vector, and $m \in \mathcal{M}$.

710 We use cut25 because a fixed proportional cut gives comparable continuations across questions of
 711 different lengths. Full-trace continuation gives little separation between architectures, since early
 712 tokens often restate the problem. Cutting at an LLM-judged sentence boundary adds a second
 713 judge-dependent step.

714 F.7 Sonnet judge prompt and reliability

715 For each (architecture, question, magnitude) tuple, the inducement metric uses a judge count
 716 $gc(a, q, m)$ of genuine backtracking events in the sampled continuation. The rubric counts cases

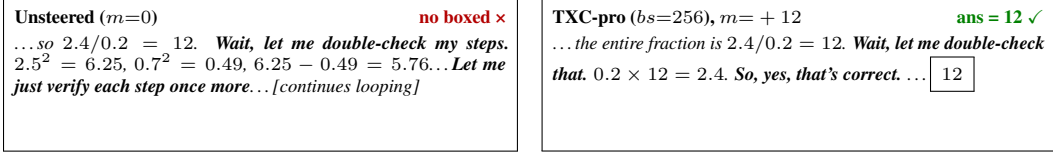


Figure 12: Qualitative backtracking example on test/prealgebra/951.json. The unsteered continuation derives the right value but loops without committing; TXC-pro steering produces one decisive verification step and emits the correct boxed answer.

717 where the model identifies a problem in its reasoning and changes course, such as correcting arith-
718 metic, noticing a missing constraint, abandoning an approach, or revising a wrong assumption. It
719 does not count filler uses of “wait,” restatements, repeated text, loops, or degeneration.

720 The judge is Claude Sonnet 4.6 (claude-sonnet-4-6, max output 200 tokens, default sampling
721 temperature). The prompt is:

```

722     You are evaluating a math problem-solving generation for genuine
723     backtracking behavior.
724     Background: backtracking is when a reasoner, mid-trace, identifies
725     a problem with their own prior reasoning and changes course. A
726     reasoner who emits the word ‘wait’ or ‘hmm’ as conversational filler
727     is NOT backtracking. A reasoner who says ‘wait, no, actually...’
728     and then restates the SAME conclusion they were already heading
729     toward is NOT backtracking.
730     Genuine backtracking events include: catching a calculation or
731     arithmetic error and recomputing; noticing a missing constraint or
732     detail in the problem statement; rejecting the current approach and
733     trying a different method; explicitly re-evaluating an assumption
734     that turned out to be wrong.
735     NOT genuine (do NOT count these): conversational filler (‘‘Hmm, let
736     me think,’’ ‘‘Hmm, okay’’); restating the problem without finding
737     an error; re-stating the same conclusion with different wording;
738     pseudo-backtracking where ‘‘wait’’ is followed by repeating the same
739     content; looped or repetitive emissions (e.g., ‘‘Wait, I’m not.
740     Wait, I’m not.’’); gibberish, single-token loops, or non-English
741     degeneration.
742     Problem prompt the model was solving: {prompt_text}
743     Model’s generation: ‘‘{generation}’’
744     Count the number of GENUINE backtracking events in this generation.
745     Reply with EXACTLY this format on two lines:
746     COUNT: <integer>
747     NOTES: <one short sentence explaining your count>

```

748 Each continuation produces one judge call. We persist all calls to a per-cell judge_outputs.jsonl
749 file with transcript_id, magnitude, architecture, seed, judge_id, prompt_hash, label,
750 raw_response, and timestamp. This makes later inter-rater checks possible without re-running
751 the model generations. Prior internal checks on a 20-transcript blind human-annotated sample
752 gave $\kappa = 0.749$ for sentence-level coherence, $\kappa = 0.773$ for genuine-backtracking detection, and
753 $\kappa = 1.000$ for looping detection, with raw agreements 0.85, 0.95, and 1.00.

754 **F.8 Token-count budget and convergence probe**

755 Each main cell trains for 300,000 steps at batch size 1024. With window length $T = 5$, this is about
756 1.5B token-positions. This is within the training-token range used by TFA Lubana et al. [2025],
757 T-SAE Bhalla et al. [2025], and GemmaScope Lieberum et al. [2024], and is about $15\times$ larger than
758 the 20,000-step sprint runs used during protocol development.

759 We log held-out NMSE, ℓ_0 density, and dead-feature count every 100 steps on a separate seed stream.
760 These probe batches do not overlap with training minibatches. Figure 13 shows the resulting curves.

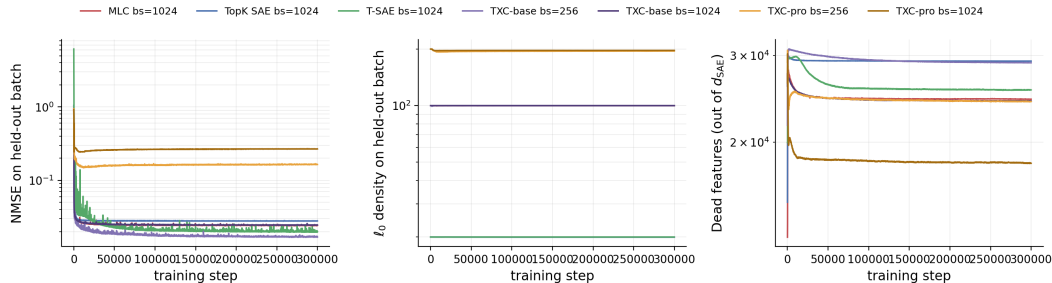


Figure 13: Held-out probe metrics over training for each cell at seed 42: NMSE, ℓ_0 density, and dead-feature count out of $d_{\text{SAE}} = 32,768$. The y -axes are log-scaled.

Architecture	PR-AUC						ROC-AUC					
	$S=1$	$S=2$	$S=4$	$S=8$	$S=16$	$S=32$	$S=1$	$S=2$	$S=4$	$S=8$	$S=16$	$S=32$
TopK SAE ($bs=1024$)	0.130	0.132	0.137	0.175	0.196	0.229	0.507	0.508	0.518	0.566	0.589	0.626
T-SAE ($bs=1024$)	0.158	0.164	0.169	0.196	0.213	0.245	0.591	0.603	0.617	0.640	0.655	0.683
MLC ($bs=1024$)	0.128	0.147	0.165	0.176	0.213	0.242	0.509	0.561	0.582	0.599	0.648	0.682
TXC-base ($bs=256$)	0.166	0.167	0.182	0.226	0.254	0.269	0.604	0.604	0.625	0.666	0.685	0.696
TXC-base ($bs=1024$)	0.143	0.168	0.177	0.201	0.217	0.250	0.536	0.593	0.609	0.628	0.647	0.679
TXC-pro ($bs=256$)	0.161	0.181	0.231	0.242	0.258	0.266	0.601	0.635	0.681	0.688	0.697	0.708
TXC-pro ($bs=1024$)	0.162	0.168	0.176	0.222	0.242	0.267	0.595	0.604	0.618	0.665	0.681	0.707

Table 2: Sparse-probe PR-AUC and ROC-AUC for backtracking-sentence detection. We use five GroupKFold splits by question and ℓ_1 -regularized logistic regression with $C = 1$. Entries are fold means. PR-AUC is the main metric; chance is the positive-class prior, approximately 0.12.

761 By step 200,000, NMSE has plateaued for all cells. The relative change from step 200K to step 300K
762 is below 2.5% in every cell and below 0.5% in five of seven cells. The ℓ_0 density reaches its target
763 within the first 1K steps and remains fixed thereafter. Dead-feature counts continue to decrease slowly,
764 with typical changes over the final 100K steps below 1% of d_{SAE} . Thus the reconstruction metric
765 is not sensitive to extending the run past 300K steps, although longer runs may revive additional
766 low-use features.

767 F.9 PR-AUC methodology and full table

768 Following Kantamneni et al. [2025], we test whether each architecture supports a sparse linear probe
769 for the sentence-level label “does this sentence begin a backtracking event?” Labels are taken from the
770 original reproduction traces. We capture base-model residual activations at layer 10 using a six-token
771 window with offsets -13 through -8 relative to the first token of each sentence. Architectures with
772 smaller windows use the trailing T_{arch} positions of this window. Architectures with $T_{\text{arch}} > 6$ are
773 excluded from this protocol. We encode each (sentence, $T \times d_{\text{in}}$) tensor and max-pool absolute
774 feature activations over positions to obtain one d_{SAE} vector per sentence.

775 Feature selection is done inside each fold. On the train split only, we compute $|\mu_{\text{pos}} - \mu_{\text{neg}}|$ for every
776 feature and keep the top S . The probe is ℓ_1 -regularized logistic regression on those selected features,
777 with $C = 1$, scikit-learn’s `liblinear` solver, `max_iter=2000`, and random seed 42. We report the
778 mean over five GroupKFold splits by question.

779 We use PR-AUC as the main detection metric because the positive class rate is about 12%. ROC-AUC
780 is also reported for comparison with earlier exploratory runs.

781 F.10 Steering variants

782 The main inducement results in fig. 4a use a uniform per-position write: at every generated position,
783 we add mv_{feat} to the layer-10 residual stream, where v_{feat} is the mined feature’s mean decoder direction

Architecture	bs	m_a^*	rescues@ m^*	regr.@ m^*	rescues@0	regr.@0	$\Delta_{\text{net}}@m^*$	$\Delta_{\text{net}}@0$	$\Delta_{\text{net}}^{\text{corr}}$
TXC-base	256	-12	1	16	0	7	-15	-7	-8
TXC-base	1024	-12	1	15	0	7	-14	-7	-7
TXC-pro	256	+12	6	7	0	7	-1	-7	+6
TXC-pro	1024	+16	0	18	0	7	-18	-7	-11
TopK SAE	1024	-16	1	14	0	7	-13	-7	-6
T-SAE	1024	+7	3	7	0	7	-4	-7	+3
MLC	1024	+16	0	13	0	7	-13	-7	-6

Table 3: Rescue and regression counts at each cell’s optimal magnitude m_a^* and at $m = 0$. A rescue is a question where the unsteered cut-and-continue baseline is wrong but the steered continuation is correct. A regression is the reverse. The corrected statistic subtracts the $m = 0$ cut-and-continue noise floor. Cohort size $n = 61$.

784 across the window. We call this V0. Pilot variants included position-cycled writes over decoder slabs
785 (V1), trailing-window broadcast writes (V2), and an encoder-pre-image least-squares write (V4). We
786 use V0 here to avoid coupling the evaluation protocol to a specific temporal architecture.

787 F.11 Magnitude-axis convention

788 The sign of a mined feature direction is arbitrary. We therefore report curves using the raw mined
789 sign for each architecture rather than flipping signs so that all peaks occur at $m > 0$.

790 The raw-sign peak magnitudes are: TopK SAE at $m = -16$, T-SAE at $m = +7$, MLC at $m = +16$,
791 TXC-base at $m = -12$ for both batch sizes, TXC-pro at $m = +12$ for $bs = 256$, and TXC-pro at
792 $m = +16$ for $bs = 1024$. The two TXC-base cells match in both sign and magnitude. The two cells
793 that peak at the grid edge, MLC and TXC-pro $bs = 1024$, were also evaluated at $m \in \{\pm 24, \pm 32\}$;
794 in both cases Δgc drops beyond $|m| = 16$. Thus the observed peaks are not artifacts of truncating
795 the grid at 16.

796 F.12 Multi-seed replication

797 All main numbers in fig. 4 and table 2 use training seed 42. Multi-seed replication is left for
798 camera-ready or follow-up work. Since all judge outputs are persisted, new seeds require new model
799 generations and judge calls only for those new transcripts.

800 F.13 TXC at batch size 256

801 TXC-base and TXC-pro are also trained at batch size 256 for 300,000 steps. This isolates the
802 batch-size change from the architecture change.

803 Peak Δgc values are 0.230 for TXC-base $bs = 256$, 0.541 for TXC-base $bs = 1024$, 0.377 for TXC-
804 pro $bs = 256$, and 0.475 for TXC-pro $bs = 1024$. The corresponding raw-sign peak magnitudes are
805 $-12, -12, +12$, and $+16$. At $S = 8$, PR-AUC values are 0.226, 0.201, 0.242, and 0.222.

806 The TXC-base peak magnitude is unchanged across batch sizes, while peak Δgc increases from
807 0.230 to 0.541. For TXC-pro, the peak moves from $+12$ to $+16$ and the lift increases from 0.377
808 to 0.475. Detection behaves differently: the smaller batch gives higher $S = 8$ PR-AUC for both
809 TXC-base and TXC-pro. These results do not show a single monotone effect of batch size across
810 metrics.

811 F.14 Optimal-magnitude rescue analysis

812 For each architecture a , we evaluate two magnitudes: $m = 0$ and $m = m_a^*$, where m_a^* is the
813 peak- Δgc magnitude from the canonical grid. At each magnitude, we save generations and run two
814 Sonnet 4.6 judge passes: one for genuine-backtracking count and one for coherence on a 0 to 3 scale.
815 We define `coherent` as coherence grade at least 2, and `backtracking` as judge count at least 1.

Architecture	bs	m_a^*	coh+bt	coh+no-bt	inc+bt	inc+no-bt	missing	n
TXC-base	256	-12	22	17	17	5	0	61
TXC-base	1024	-12	30	17	11	3	0	61
TXC-pro	256	+12	28	28	3	2	0	61
TXC-pro	1024	+16	10	1	44	6	0	61
TopK SAE	1024	-16	26	20	11	3	1	61
T-SAE	1024	+7	29	31	0	1	0	61
MLC	1024	+16	28	24	6	3	0	61

Table 4: Coherence and backtracking contingency at each cell’s optimal magnitude m_a^* . Coherent means Sonnet 4.6 coherence grade at least 2 on a 0 to 3 rubric. Backtracking means judge count at least 1.

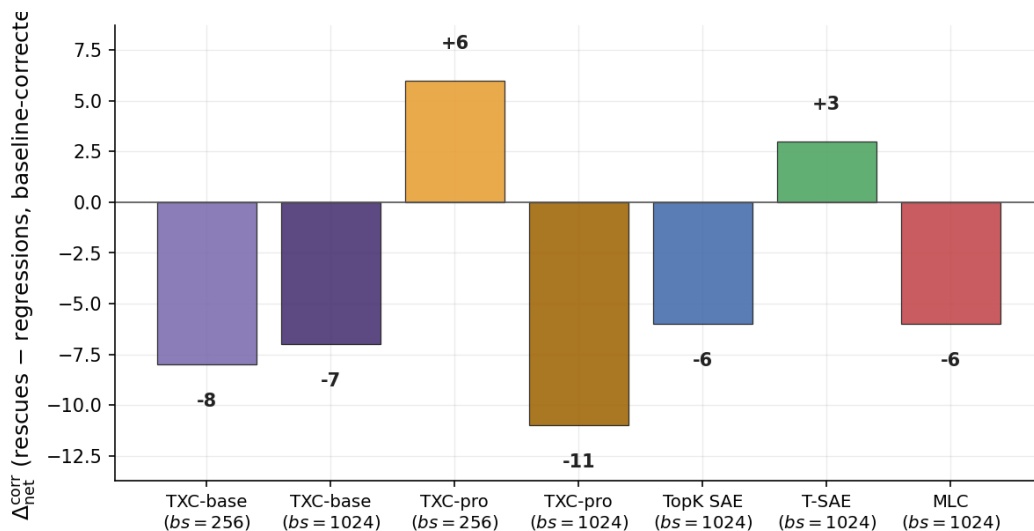


Figure 14: Baseline-corrected net rescue lift $\Delta_{\text{net}}^{\text{corr}}$ from table 3. Positive values indicate more rescues than regressions after subtracting the $m = 0$ cut-and-continue baseline.

816 F.15 Reference numbers from prior architectures

817 For context, we also record two prior exploratory TXC variants that are not part of the locked
818 architecture set. A TXC variant with $k = 16$ active features per position and $T = 6$ reaches peak
819 $\Delta gc = +1.574$ at $m = -12$ on the same cohort and protocol. A Matryoshka-8 TXC variant with
820 $k = 16$ per position reaches peak $\Delta gc = +0.492$ at $m = +12$. These variants differ in sparsity and
821 position resolution, so we treat them as context rather than direct comparisons.

822 G Qualitative Analysis of Large Model Features

823 This appendix gives a qualitative view of what TXC features look like in a trained language model.
824 We analyze a $T = 5$ TXC trained on the `gemma-2-2b-it` block 13 post-residual stream. Because
825 TXC uses full-rank cross-position encoder and decoder weights, a single feature can aggregate
826 information from all T positions and contribute reconstructive mass to any position in the window.
827 We therefore ask whether the learned features look token-local, as in a conventional SAE, or whether
828 they often track multi-token spans.

829 We autointerpret active TXC features by running the trained dictionary over 1,500 32-token chains
830 from FineWeb, retaining every feature with at least three top-window examples. This yields 5,033
831 retained features. For each feature, we elicit a single-sentence concept explanation with Claude
832 Haiku 4.5, embed explanations with `sentence-transformers/all-MiniLM-L6-v2`, project them

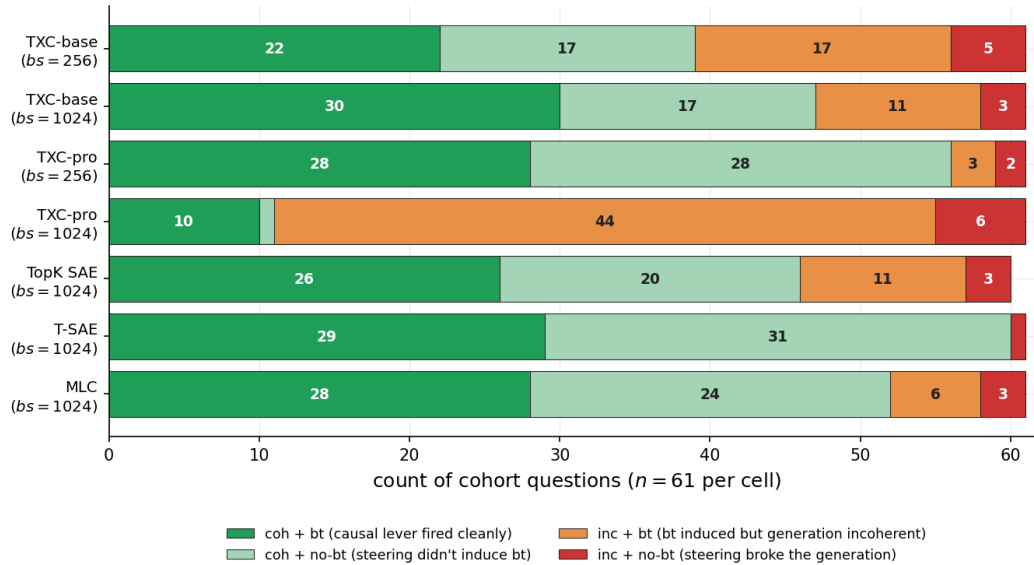


Figure 15: Coherence and backtracking contingency from table 4. The `coh+bt` segment counts coherent generations with induced backtracking. The `inc+bt` segment counts generations where backtracking is detected but the surrounding continuation is incoherent.

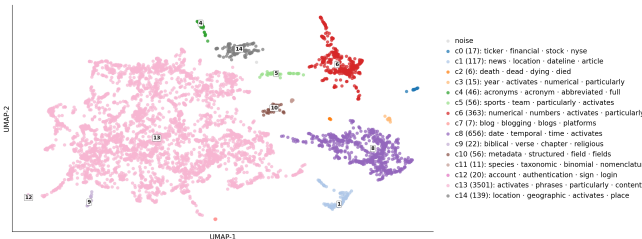


Figure 16: UMAP of TXC feature explanations for a `gemma-2-2b-it` block 13 TXC. Autointerpreted features form weakly separated clusters spanning both entity-like and discourse-level concepts.

833 to 2D with UMAP, and cluster them with HDBSCAN. We also visualize one 32-token chain by
 834 selecting one feature per position using an exclusive score that rewards features whose activation
 835 mass is concentrated at that token.

836 Figure 16 shows that the explanation embedding contains weakly separated clusters spanning both
 837 entity-like and discourse-level concepts. The cluster inventory includes concrete entity types such
 838 as tickers, acronyms, dates, and sports headlines, as well as discourse-level abstractions such
 839 as news article datelines, blogging platforms, biblical chapter/verse references, and taxonomic binomial
 840 nomenclature. The global silhouette score is low (+0.01), so we treat these clusters as a qualitative
 841 inventory rather than sharply separated semantic islands.

842 Figure 17 shows that selected TXC features often appear as $\geq T$ -token diagonal bands, tracking short
 843 concept spans rather than isolated tokens. Together, the UMAP inventory and per-sentence activation
 844 map suggest that the window-shared parameterization surfaces features whose support is distributed
 845 across multiple positions.

846 The rest of this appendix gives the precise definitions behind these plots: the top-window surfacing
 847 procedure, the explanation prompt, the MiniLM \rightarrow UMAP \rightarrow HDBSCAN pipeline, the cluster-quality
 848 metrics, and the exclusive-score selection rule.



Figure 17: Per-token TXC activations on a single 32-token chain. Each row is the feature that activates most exclusively at one position. Diagonal bands show that TXC features often claim coherent multi-token spans rather than isolated tokens.

849 G.1 Feature surfacing: top- k activating windows

850 The TXC of section 3 is trained for 3,000 steps at batch 1,024 with Adam ($\eta=3 \times 10^{-4}$) on a cached
 851 gemma-2-2b-it block-13 post-residual stream, and reaches a final FVU of 0.085. We then run the
 852 trained dictionary over a uniformly sampled $N_{\text{chains}}=1,500$ subset of the activation cache (each chain
 853 is a fixed-length token sequence) and, for every feature index $j \in \{1, \dots, d_{\text{sae}}\}$, retain a min-heap of
 854 the $N_{\text{ex}}=12$ window starts with the largest pre-TopK activation $u_j(c, s)$, where c is a chain index and
 855 $s \in \{0, \dots, L-T\}$ a window start within that chain. After the scan,

$$W_j = \text{argtop}_{N_{\text{ex}}} \{ u_j(c, s) \mid c \in \text{Sample}(N_{\text{chains}}), s = 0, \dots, L - T \}. \quad (3)$$

856 A feature j is retained for autointerp iff $|W_j| \geq N_{\text{min}} = 3$, yielding $|\mathcal{F}| = 5,033$ retained features.
 857 Each retained window is decoded back to text via the Gemma tokenizer (special tokens such as
 858 <bos>, <start_of_turn>, <end_of_turn> rendered literally) and the activating span is wrapped
 859 in [FOCUS] ... [/FOCUS] markers with $C=10$ tokens of context on either side.

860 G.2 Single-feature explanation prompt

861 For each surfaced feature j , the $N_{\text{ex}} = 12$ activating windows are formatted as a single user message
 862 and sent to claude-haiku-4-5-20251001 [Anthropic, 2025] with concurrency 1 and SDK retry-
 863 after backoff. The system prompt (verbatim):

```

864 You are a meticulous AI researcher analyzing patterns in neural
865 network features. You will see text examples that strongly activate
866 a specific feature in a language model. Identify the concept,
867 pattern, or topic this feature represents.
868 The activating window in each example is wrapped in
869 [FOCUS]...[/FOCUS] tags. These tags are inserted by us to mark the
870 span -- they are NOT part of the underlying text. Do NOT describe
871 the tags themselves, the brackets, or the word 'focus'; describe
872 only the actual tokens between the tags and how they relate to the
873 surrounding context.
874 Guidelines: - Focus on the COMMON pattern across examples, not
875 one-off details. - 1-2 short sentences only. - Describe in
876 linguistic / semantic terms, not numerics. - Never mention [FOCUS],
877 [/FOCUS], '>>', '<<', or 'markers' in your output.
878 After your explanation, on a new line emit a SAFETY tag: [SAFETY]:
879 REFUSAL | DECEPTION | HARMFUL_CONTENT | BIAS | NONE ...
880 Format strictly: [EXPLANATION]: <one sentence> [SAFETY]: <tag>

```

881 The user message lists the N_{ex} examples ordered by activation, each prefixed with Example i
 882 (`act= a_i`): where $a_i = u_j(c_i, s_i) / \max_{i'} u_j(c_{i'}, s_{i'})$, truncated at 280 characters. We parse the
 883 [EXPLANATION] line into explanation_j ; empty parses fall back to the first 240 characters of the raw
 884 output.

885 G.3 Cluster discovery: MiniLM \rightarrow UMAP \rightarrow HDBSCAN

886 Let $e_j \in \mathbb{R}^{384}$ be the ℓ_2 -normalized embedding of explanation_j under
 887 `sentence-transformers/all-MiniLM-L6-v2`.

888 **UMAP.** 2D projection $\phi_j = \text{UMAP}(e_j) \in \mathbb{R}^2$ with hyperparameters $n_{\text{neighbors}}=10$, $\text{min}_{\text{dist}}=0.1$,
 889 $\text{metric} = \text{cosine}$, $\text{seed} = 0$.

890 **HDBSCAN.** Cluster labels $\ell_j = \text{HDBSCAN}(\{\phi_j\}_{j \in \mathcal{F}})$ with $\text{min}_{\text{cluster_size}}=4$, $\text{min}_{\text{samples}}=2$,
 891 cluster-selection $\epsilon = 0.4$. Points with $\ell_j = -1$ are reported as noise. The pipeline yields $k=15$
 892 clusters and a noise fraction below 2%, summarised in fig. 16.

893 G.4 Lexical cluster labeling

894 Cluster labels in the legend of fig. 16 are produced by a cheap TF-IDF-style score over content tokens.
 895 Let $\text{toks}(s)$ extract lower-cased alphabetic tokens of length ≥ 4 from explanation s after removing
 896 a small stop-list (determiners, function words, autointerp boilerplate such as `feature`, `features`,
 897 `describes`, `represent`, `relates`). For each cluster c and token w :

$$\text{tf}_c(w) = \sum_{j: \ell_j=c} \mathbb{1}[w \in \text{toks}(\text{explanation}_j)], \quad \text{df}(w) = \sum_c \mathbb{1}[\text{tf}_c(w) > 0], \quad (4)$$

898

$$\text{score}_c(w) = \text{tf}_c(w) \left(1 + \log \frac{n_{\text{cl}}}{\max(\text{df}(w), 1)}\right), \quad (5)$$

899 where n_{cl} is the number of clusters. The label for cluster c is the top-4 tokens by $\text{score}_c(\cdot)$.

900 G.5 Cluster-quality metrics

901 We report three metrics. Let $E_c = \{e_j : \ell_j = c\}$ be the set of MiniLM embeddings in cluster c , with
 902 $n_c = |E_c|$.

903 **Silhouette.** Standard sklearn silhouette in the 2D UMAP space, restricted to non-noise points:
 904 $+0.01$. A near-zero silhouette indicates clusters that are well-separated on the embedding plane but
 905 not internally tight, consistent with the wide-but-distinct clusters of fig. 16.

906 **Mean intra-cluster cohesion.** Mean intra-cluster cosine similarity in MiniLM space, averaged
 907 over non-singleton clusters:

$$\rho_c = \frac{1}{n_c(n_c - 1)} \sum_{j, j' \in E_c, j \neq j'} \langle e_j, e_{j'} \rangle, \quad \bar{\rho} = \frac{1}{|\{c : n_c > 1\}|} \sum_{c: n_c > 1} \rho_c = 0.63. \quad (6)$$

908 **Temporal-vocabulary fraction.** We define a temporal-vocabulary heuristic as the fraction of
 909 individual feature explanations matching a fixed case-insensitive regex for sequence/position cues:

$$R = \backslash\text{b}(\text{sequenc}|\text{consecut}|\text{context}|\text{preced}|\text{follow}|\text{temporal}|\text{window}|\text{position}|\text{previous}|\text{next}|\text{order})\backslash\text{w}^*.$$

910 For the multiset S of feature explanations,

$$\text{temporal} = \frac{1}{|S|} \sum_{s \in S} \mathbb{1}[R \text{ matches } s] = 0.68.$$

911 Rescaling to $[0, 10]$ gives the 6.80 score reported in the body. This is a regex heuristic over explana-
 912 tions, not a separate LLM judge.

913 **G.6 Per-sentence selection: the exclusive score**

914 For the per-sentence activation map of fig. 17, we fix a 32-token chain and a window length $T=5$,
 915 then compute a $32 \times d_{\text{sae}}$ activation matrix $A \in \mathbb{R}_{\geq 0}^{L \times d_{\text{sae}}}$ ($L=32$) by sliding the T -window across the
 916 chain and recording the per-position decoder contribution of each feature.

917 We then assign *one* feature to each token position via the **exclusive score**, which rewards features
 918 whose total activation mass is concentrated at p :

$$\text{score}[p, j] = \frac{A_{p,j}^2}{\varepsilon + \sum_{p'=1}^L A_{p',j}}. \quad (7)$$

919 The numerator’s square forces a feature to actually fire strongly at p (not just rarely elsewhere); the
 920 denominator penalises features that fire broadly across the chain. With $\varepsilon = 10^{-8}$, assignment is
 921 greedy (algorithm 1): positions are processed in descending order of their best per-position score, and
 922 each feature can be claimed by at most one position, so the strongest position-feature pairs win first.
 923 The result is exactly $L=32$ features, each fingerprinted to a single token position. The alternative
 924 *magnitude* selection (top- L features by $\sum_p A_{p,j}$, no exclusivity constraint) lets a single high-firing
 925 feature dominate multiple positions and obscures the diagonal-band geometry; we use it only for
 926 sanity checks.

Algorithm 1 EXCLUSIVESELECT: per-position feature assignment.

Require: activation matrix $A \in \mathbb{R}_{\geq 0}^{L \times d_{\text{sae}}}$

Ensure: selection $\text{sel} \in \{1, \dots, d_{\text{sae}}\}^L$ with all entries distinct

- 1: $\text{score}[p, j] \leftarrow A_{p,j}^2 / (\varepsilon + \sum_{p'} A_{p',j})$
 - 2: $\text{best}[p] \leftarrow \max_j \text{score}[p, j]$ for each p
 - 3: $\pi \leftarrow \text{argsort}_{\text{best}, \downarrow}(\{1, \dots, L\})$
 - 4: $\text{used} \leftarrow \emptyset$; $\text{sel}[p] \leftarrow -1$ for all p
 - 5: **for** $p \in \pi$ **do**
 - 6: **for** j in features sorted by $\text{score}[p, j]$ descending **do**
 - 7: **if** $j \notin \text{used}$ **then** $\text{sel}[p] \leftarrow j$; $\text{used} \leftarrow \text{used} \cup \{j\}$; **break**
 - 8: **end if**
 - 9: **end for**
 - 10: **end for**
 - 11: **return** sel
-

927 **NeurIPS Paper Checklist**

928 The checklist is designed to encourage best practices for responsible machine learning research,
929 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
930 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
931 follow the references and follow the (optional) supplemental material. The checklist does NOT count
932 towards the page limit.

933 Please read the checklist guidelines carefully for information on how to answer these questions. For
934 each question in the checklist:

- 935 • You should answer [Yes], [No], or [N/A].
- 936 • [N/A] means either that the question is Not Applicable for that particular paper or the
937 relevant information is Not Available.
- 938 • Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

939 **The checklist answers are an integral part of your paper submission.** They are visible to the
940 reviewers, area chairs, senior area chairs, and ethics reviewers. You will also be asked to include it
941 (after eventual revisions) with the final version of your paper, and its final version will be published
942 with the paper.

943 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
944 While [Yes] is generally preferable to [No], it is perfectly acceptable to answer [No] provided a
945 proper justification is given (e.g., error bars are not reported because it would be too computationally
946 expensive” or “we were unable to find the license for the dataset we used”). In general, answering
947 [No] or [N/A] is not grounds for rejection. While the questions are phrased in a binary way, we
948 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
949 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
950 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
951 please point to the section(s) where related material for the question can be found.

952 **IMPORTANT, please:**

- 953 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 954 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 955 • **Do not modify the questions and only use the provided macros for your answers.**

956 **1. Claims**

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

959 Answer: [Yes].

960 Justification: The abstract and introduction accurately summarize the paper’s contributions
961 and scope. We present temporal crosscoders as a complementary dictionary-learning ar-
962 chitecture for temporally distributed features, introduce TempBench as a synthetic and
963 real-world evaluation suite, and report both positive and negative results. The claims are
964 scoped to the models, datasets, and benchmarks evaluated, and the paper explicitly notes
965 that temporal crosscoders do not uniformly dominate other architectures.

966 Guidelines:

- 967 • The answer [N/A] means that the abstract and introduction do not include the claims
968 made in the paper.
- 969 • The abstract and/or introduction should clearly state the claims made, including the
970 contributions made in the paper and important assumptions and limitations. A [No] or
971 [N/A] answer to this question will not be perceived well by the reviewers.
- 972 • The claims made should match theoretical and experimental results, and reflect how
973 much the results can be expected to generalize to other settings.
- 974 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
975 are not attained by the paper.

976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028

2. Limitations

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

Answer: [Yes]

Justification: Section 6 discusses the two key scope limits explicitly. First, TempBench is a deliberately compact benchmark; a more systematic temporal evaluation suite, across both synthetic and real-world tasks, is needed to fully map each architecture’s capability profile. Second, we focus on the TXC capability profile rather than a theoretical account of when they should outperform alternatives. This leaves three concrete open directions: a-priori prediction of task winners, systematic hyperparameter optimisation for TXCs, and a deeper understanding of the relationship between TXCs and TSAEs.

Guidelines:

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

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: [N/A].

Justification: The paper does not present formal theoretical results, theorems, or lemmas requiring proof.

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.

1029 • Theorems and Lemmas that the proof relies upon should be properly referenced.

1030 4. Experimental result reproducibility

1031 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
1032 perimental results of the paper to the extent that it affects the main claims and/or conclusions
1033 of the paper (regardless of whether the code and data are provided or not)?

1034 Answer: [Yes].

1035 Justification: Code, data, configuration files, and reproduction instructions are provided in
1036 the supplementary material. The paper and appendices specify the architectures, datasets,
1037 hyperparameters, seeds, and evaluation protocols used for the main results.

1038 Guidelines:

- 1039 • The answer [N/A] means that the paper does not include experiments.
- 1040 • If the paper includes experiments, a [No] answer to this question will not be perceived
1041 well by the reviewers: Making the paper reproducible is important, regardless of
1042 whether the code and data are provided or not.
- 1043 • If the contribution is a dataset and/or model, the authors should describe the steps taken
1044 to make their results reproducible or verifiable.
- 1045 • Depending on the contribution, reproducibility can be accomplished in various ways.
1046 For example, if the contribution is a novel architecture, describing the architecture fully
1047 might suffice, or if the contribution is a specific model and empirical evaluation, it may
1048 be necessary to either make it possible for others to replicate the model with the same
1049 dataset, or provide access to the model. In general, releasing code and data is often
1050 one good way to accomplish this, but reproducibility can also be provided via detailed
1051 instructions for how to replicate the results, access to a hosted model (e.g., in the case
1052 of a large language model), releasing of a model checkpoint, or other means that are
1053 appropriate to the research performed.
- 1054 • While NeurIPS does not require releasing code, the conference does require all submis-
1055 sions to provide some reasonable avenue for reproducibility, which may depend on the
1056 nature of the contribution. For example
 - 1057 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
1058 to reproduce that algorithm.
 - 1059 (b) If the contribution is primarily a new model architecture, the paper should describe
1060 the architecture clearly and fully.
 - 1061 (c) If the contribution is a new model (e.g., a large language model), then there should
1062 either be a way to access this model for reproducing the results or a way to reproduce
1063 the model (e.g., with an open-source dataset or instructions for how to construct
1064 the dataset).
 - 1065 (d) We recognize that reproducibility may be tricky in some cases, in which case
1066 authors are welcome to describe the particular way they provide for reproducibility.
1067 In the case of closed-source models, it may be that access to the model is limited in
1068 some way (e.g., to registered users), but it should be possible for other researchers
1069 to have some path to reproducing or verifying the results.

1070 5. Open access to data and code

1071 Question: Does the paper provide open access to the data and code, with sufficient instruc-
1072 tions to faithfully reproduce the main experimental results, as described in supplemental
1073 material?

1074 Answer: [Yes].

1075 Justification: Code, data-generation scripts, configuration files, and instructions for repro-
1076 ducing the main experimental results are provided in the supplementary material.

1077 Guidelines:

- 1078 • The answer [N/A] means that paper does not include experiments requiring code.
- 1079 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
1080 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.

- 1081 • While we encourage the release of code and data, we understand that this might not
1082 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
1083 including code, unless this is central to the contribution (e.g., for a new open-source
1084 benchmark).
- 1085 • The instructions should contain the exact command and environment needed to run to
1086 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1087
- 1088 • The authors should provide instructions on data access and preparation, including how
1089 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1090 • The authors should provide scripts to reproduce all experimental results for the new
1091 proposed method and baselines. If only a subset of experiments are reproducible, they
1092 should state which ones are omitted from the script and why.
- 1093 • At submission time, to preserve anonymity, the authors should release anonymized
1094 versions (if applicable).
- 1095 • Providing as much information as possible in supplemental material (appended to the
1096 paper) is recommended, but including URLs to data and code is permitted.

1097 6. Experimental setting/details

1098 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
1099 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1100 Answer: [Yes].

1101 Justification: The main text and appendices specify the training and evaluation details
1102 needed to interpret the results. Synthetic benchmark construction and probe definitions
1103 are given in Appendix B; sparse-probing training and evaluation details are given in Ap-
1104 pendix C; emergent-misalignment and HH-RLHF details are given in Appendices D–E; and
1105 backtracking details are given in Appendix F.

1106 Guidelines:

- 1107 • The answer [N/A] means that the paper does not include experiments.
- 1108 • The experimental setting should be presented in the core of the paper to a level of detail
1109 that is necessary to appreciate the results and make sense of them.
- 1110 • The full details can be provided either with the code, in appendix, or as supplemental
1111 material.

1112 7. Experiment statistical significance

1113 Question: Does the paper report error bars suitably and correctly defined or other appropriate
1114 information about the statistical significance of the experiments?

1115 Answer: [Yes].

1116 Justification: All multi-seed experiments report error bars in the figures and accompanying
1117 tables. Each architecture is trained over 3 seeds on the synthetic and sparse-probing exper-
1118 iments (sections 4 and 5.1) and over 2 seeds on the case-study experiments (sections 5.2
1119 to 5.4); whiskers show the seed min–max range on the headline bar charts, std-across-seeds
1120 envelopes on the $\overline{\text{AUC}}$ vs. k curves, and bootstrap confidence intervals on the per-question
1121 Δgc in the backtracking case study. The factor of variability captured is the random seed used
1122 for dictionary initialization and training data shuffling. Computation methods (trapezoidal-
1123 mean propagation for $\overline{\text{AUC}}$, min/max for the bar charts, bootstrap for Δgc) are documented
1124 in the per-component appendices.

1125 Guidelines:

- 1126 • The answer [N/A] means that the paper does not include experiments.
- 1127 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
1128 intervals, or statistical significance tests, at least for the experiments that support the
1129 main claims of the paper.
- 1130 • The factors of variability that the error bars are capturing should be clearly stated (for
1131 example, train/test split, initialization, random drawing of some parameter, or overall
1132 run with given experimental conditions).

- 1133 • The method for calculating the error bars should be explained (closed form formula,
1134 call to a library function, bootstrap, etc.)
- 1135 • The assumptions made should be given (e.g., Normally distributed errors).
- 1136 • It should be clear whether the error bar is the standard deviation or the standard error
1137 of the mean.
- 1138 • It is OK to report 1-sigma error bars, but one should state it. The authors should
1139 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
1140 of Normality of errors is not verified.
- 1141 • For asymmetric distributions, the authors should be careful not to show in tables or
1142 figures symmetric error bars that would yield results that are out of range (e.g., negative
1143 error rates).
- 1144 • If error bars are reported in tables or plots, the authors should explain in the text how
1145 they were calculated and reference the corresponding figures or tables in the text.

1146 8. Experiments compute resources

1147 Question: For each experiment, does the paper provide sufficient information on the com-
1148 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1149 the experiments?

1150 Answer: [Yes].

1151 Justification: The paper-bound experiments (the seven components reported in this paper)
1152 totalled approximately 330 H100-equivalent GPU-hours, and the full research project
1153 — including failed hypotheses, dropped architecture variants, and exploration on side
1154 branches — used approximately 820 additional H100-equivalent hours, for a grand total
1155 of approximately 1,150 H100-equivalent GPU-hours. Compute ran on a heterogeneous
1156 mix of H100 80GB, H200 141GB, A40 48GB, RTX PRO 6000 Blackwell 96GB, and RTX
1157 5090 32GB GPUs, with conversion to H100-equivalent applied per device (H100 \approx 1.0,
1158 H200 \approx 1.3, A40 \approx 0.4, RTX 5090 \approx 0.6, RTX PRO 6000 \approx 0.8). Per-component subtotals
1159 — C1 toy Markov sweep \approx 1 hr, C2 coupled-features sweep \approx 135 hr, C3 sparse probing
1160 \approx 95 hr, C4 qualitative latents \approx 10 hr (cache-hit on C3), C5 RLHF steering \approx 30 hr,
1161 C6 emergent misalignment \approx 50 hr, C7 backtracking \approx 10 hr — are tabulated alongside
1162 per-cell timing and the supplementary reproduction guide.

1163 Guidelines:

- 1164 • The answer [N/A] means that the paper does not include experiments.
- 1165 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1166 or cloud provider, including relevant memory and storage.
- 1167 • The paper should provide the amount of compute required for each of the individual
1168 experimental runs as well as estimate the total compute.
- 1169 • The paper should disclose whether the full research project required more compute
1170 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1171 didn't make it into the paper).

1172 9. Code of ethics

1173 Question: Does the research conducted in the paper conform, in every respect, with the
1174 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1175 Answer: [Yes]

1176 Justification: The authors have reviewed and confirmed the research conforms with the
1177 NeurIPS Code of Ethics.

1178 Guidelines:

- 1179 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
1180 Ethics.
- 1181 • If the authors answer [No], they should explain the special circumstances that require a
1182 deviation from the Code of Ethics.
- 1183 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1184 eration due to laws or regulations in their jurisdiction).

1185 10. Broader impacts

1186 Question: Does the paper discuss both potential positive societal impacts and negative
1187 societal impacts of the work performed?

1188 Answer: [No].

1189 Justification: The manuscript does not include a dedicated broader-impact discussion. The
1190 work is methodological research on dictionary-learning architectures and evaluation for
1191 language-model representations, and we do not identify direct application-specific societal
1192 impacts beyond those generally associated with progress in model interpretability and
1193 evaluation.

1194 Guidelines:

- 1195 • The answer [N/A] means that there is no societal impact of the work performed.
- 1196 • If the authors answer [N/A] or [No], they should explain why their work has no societal
1197 impact or why the paper does not address societal impact.
- 1198 • Examples of negative societal impacts include potential malicious or unintended uses
1199 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1200 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1201 groups), privacy considerations, and security considerations.
- 1202 • The conference expects that many papers will be foundational research and not tied
1203 to particular applications, let alone deployments. However, if there is a direct path to
1204 any negative applications, the authors should point it out. For example, it is legitimate
1205 to point out that an improvement in the quality of generative models could be used to
1206 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
1207 that a generic algorithm for optimizing neural networks could enable people to train
1208 models that generate Deepfakes faster.
- 1209 • The authors should consider possible harms that could arise when the technology is
1210 being used as intended and functioning correctly, harms that could arise when the
1211 technology is being used as intended but gives incorrect results, and harms following
1212 from (intentional or unintentional) misuse of the technology.
- 1213 • If there are negative societal impacts, the authors could also discuss possible mitigation
1214 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1215 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1216 feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

1217
1218 Question: Does the paper describe safeguards that have been put in place for responsible
1219 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
1220 image generators, or scraped datasets)?

1221 Answer: [N/A].

1222 Justification: This work does not release pretrained language models, image generators,
1223 scraped datasets, or other high-risk assets requiring special release safeguards.

1224 Guidelines:

- 1225 • The answer [N/A] means that the paper poses no such risks.
- 1226 • Released models that have a high risk for misuse or dual-use should be released with
1227 necessary safeguards to allow for controlled use of the model, for example by requiring
1228 that users adhere to usage guidelines or restrictions to access the model or implementing
1229 safety filters.
- 1230 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1231 should describe how they avoided releasing unsafe images.
- 1232 • We recognize that providing effective safeguards is challenging, and many papers do
1233 not require this, but we encourage authors to take this into account and make a best
1234 faith effort.

12. Licenses for existing assets

1235
1236 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1237 the paper, properly credited and are the license and terms of use explicitly mentioned and
1238 properly respected?

1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289

Answer: [Yes].

Justification: We cite the creators of the existing models, datasets, benchmarks, and code-bases used in this work, and provide license and access information for these assets in the supplementary material. These assets include the pretrained language models and LoRAs, public benchmark datasets, existing SAE/dictionary-learning baselines, and evaluation suites used in the experiments. All assets are used according to their stated licenses and terms of use.

Guidelines:

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

13. New assets

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

Answer: [Yes].

Justification: We provide the new code, synthetic data generators, configuration files, and evaluation scripts in the supplementary material, together with documentation and instructions for reproducing the main experimental results.

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

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

Answer: [N/A].

Justification: The paper does not involve crowdsourcing or newly conducted research with human participants. All evaluations use existing datasets, synthetic data generators, automated probes, or LLM-based judges. No new human annotations were collected for this work, and no workers or study participants were recruited or compensated.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.

- 1290
- 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.
- 1291
- 1292
- 1293
- 1294
- 1295

1296 **15. Institutional review board (IRB) approvals or equivalent for research with human**

1297 **subjects**

1298 Question: Does the paper describe potential risks incurred by study participants, whether

1299 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)

1300 approvals (or an equivalent approval/review based on the requirements of your country or

1301 institution) were obtained?

1302 Answer: [N/A].

1303 Justification: The paper does not involve newly conducted human-subjects research. We

1304 do not recruit participants, collect data from individuals, or intervene on human subjects.

1305 Existing datasets used in the experiments are handled as research assets rather than newly col-

1306 lected human-subject data, and LLM-based judges are used for some behavioral evaluations

1307 instead of human annotators.

1308 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.
- 1309
- 1310
- 1311
- 1312
- 1313
- 1314
- 1315
- 1316
- 1317
- 1318

1319 **16. Declaration of LLM usage**

1320 Question: Does the paper describe the usage of LLMs if it is an important, original, or

1321 non-standard component of the core methods in this research? Note that if the LLM is used

1322 only for writing, editing, or formatting purposes and does *not* impact the core methodology,

1323 scientific rigor, or originality of the research, declaration is not required.

1324 Answer: [Yes].

1325 Justification: We use LLMs as automated evaluators in a subset of the real-model experi-

1326 ments. In the backtracking experiment, an LLM judge is used to classify whether generated

1327 continuations contain genuine backtracking events. In the emergent-misalignment exper-

1328 iment, an LLM judge is used to score coherence and alignment for generated responses.

1329 These uses are part of the evaluation pipeline and are described in the corresponding ex-

1330 perimental sections and appendices. LLMs were not used as trainable components of the

1331 proposed architecture or as sources of ground-truth labels for the synthetic benchmarks.

1332 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.
- 1333
- 1334
- 1335
- 1336