

---

# Static Unit-Scale Bias Steering Transfers Poorly to a Reasoning-Distilled LLM

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Activation steering is a leading technique for controlling LLM behavior, but its  
2       reliability in reasoning-distilled checkpoints is unclear. We study a common  
3       intervention class—static continuous addition of a pre-generation linear-probe  
4       direction—and find that a cognitive-bias direction that controls an instruction-  
5       tuned model does not transfer cleanly to a reasoning-distilled sibling under the  
6       same protocol. We construct a 470-item contrastive benchmark spanning 11 bias  
7       categories (base-rate neglect, conjunction fallacy, framing, and others) and com-  
8       pare matched-architecture pairs (Llama-3.1-8B-Instruct vs. R1-Distill-Llama-8B,  
9       OLMo-3 7B/32B Instruct vs. Think) and Qwen-3-8B’s thinking toggle. Behavioral  
10      lure rates are scorer- and benchmark-scope dependent: under final-answer rescoring  
11      of preserved 470-item raw-response artifacts, R1-Distill has lower overall lure than  
12      Llama (25.5% vs. 33.2%) but remains highly vulnerable on base-rate and conjunc-  
13      tion items; on the full 470-item OLMo-32B scale run, lure falls from 19.6% to 0.4%.  
14      Lure suppression is not treated as correctness; accuracy and other-response rates  
15      are analyzed jointly. The main result characterizing the scoped dissociation is that  
16      probe-direction steering on the three vulnerable categories produces a monotonic  
17      37.5pp dose-response in Llama (lure rates span 31% at  $\alpha=+5$  to 69% at  $\alpha=-5$ , zero  
18      incoherent outputs), while the original static continuous intervention in R1-Distill,  
19      applied across four candidate layers including its probe peak, yields only small  
20      non-monotonic lure-rate fluctuations (5.0pp full-sweep range at L31), not a stable  
21      dose-response; uncalibrated final-answer-span P0/T0 diagnostics likewise show  
22      no endpoint effect but have near-zero prompt-prefill KL. Appendix diagnostics  
23      additionally show that, in the available OLMo-family 32B comparison, the larger  
24      Instruct checkpoint has higher scored lure rate than 7B (14.9%→19.6%) while  
25      the Think checkpoint remains near-zero (0.4%). Diagnostic analyses show that  
26      Qwen-3-8B’s hard think/no-think template induces non-transferring P0 geometries  
27      and that within-CoT linear separability is non-stationary, but these diagnostics  
28      are not treated as evidence for a template-invariant semantic axis or causal CoT  
29      stages. These results support a narrow claim: a static unit-scale probe direction  
30      that steers Llama can fail to provide comparable behavioral control in R1-Distill-  
31      Llama-8B even when the same bias distinction remains linearly decodable. They  
32      do not rule out calibrated dynamic, broader multi-layer, SAE-feature, or nonlinear  
33      interventions.

## 34 1 Introduction

35      Kahneman’s dual-process framework distinguishes fast, heuristic processing (“System 1”) from  
36      slow, deliberative processing (“System 2”) [Kahneman, 2011]. While the binary framing is an  
37      acknowledged oversimplification (modern accounts treat this as a *continuum* of processing intensity

38 [Melnikoff and Bargh, 2018, Evans and Stanovich, 2013, De Neys, 2023]), the core observation  
39 that cognitive systems flexibly allocate effort, with measurable consequences for bias susceptibility,  
40 remains robust. De Neys’ conflict detection paradigm [De Neys, 2012, 2017] demonstrates that  
41 humans often *detect* conflicts between heuristic and normative responses even when they fail to  
42 override the heuristic, a dissociation between conflict *monitoring* and conflict *resolution* with clear  
43 neural correlates in the anterior cingulate cortex [Botvinick et al., 2001].

44 LLMs exhibit analogous patterns: instruction-tuned models reproduce human-like errors on conflict  
45 items while succeeding on matched controls [Hagendorff et al., 2023], and reasoning-trained models  
46 reduce many such lures [DeepSeek-AI, 2025], though this reduction is bias-specific rather than  
47 universal [Kim et al., 2025]. A growing literature documents these parallels [Coda-Forno et al., 2025,  
48 Ziabari et al., 2025], and Huang et al. [2026] showed that bias-relevant directions are probeable  
49 and steerable. We ask a different question: what does reasoning post-training change about the  
50 *decodability*, layer location, and *behavioral writability* of those directions?

51 Yet behavior alone cannot distinguish genuine processing-mode differences from surface mimicry:  
52 chain-of-thought traces are often unfaithful, with only  $\sim 2.3\%$  of reasoning steps causally influencing  
53 the final answer [Zhao et al., 2025, Lanham et al., 2023, Boppana et al., 2026]. The critical question  
54 is whether fast/slow modes have *mechanistic correlates* in internal representations, and how those  
55 correlates change when a model is trained to reason.

56 Recent work on activation steering [Turner et al., 2023, Panickssery et al., 2024, Cox et al., 2026] has  
57 shown that targeted interventions on internal representations can modulate behavior, with pre-CoT  
58 probes achieving  $>0.9$  AUC and answer-flipping rates above 50% on instruction-tuned models [Cox  
59 et al., 2026]. This raises a natural question for the dual-process setting: if linear probes identify a  
60 direction in activation space that separates S1-like from S2-like processing, does a standard static  
61 intervention along that direction modulate bias susceptibility? If so, the representational findings  
62 move from correlational to intervention-based evidence, providing evidence that the probe direction  
63 captures more than an epiphenomenal surface-feature correlation under the tested protocol.

64 We address this question with a natural experiment: Llama-3.1-8B-Instruct and DeepSeek-R1-Distill-  
65 Llama-8B share identical architectures (32 layers, 32 attention heads, 4096 hidden dimensions) but  
66 differ in reasoning distillation training. Representational differences are therefore less likely to be  
67 driven by architecture alone. This controls architecture, not the full intervention: distillation data,  
68 CoT formatting, optimization, and answer style also change. We make six contributions:

- 69 1. A **contrastive cognitive bias benchmark** of 470 items (235 matched conflict/control pairs) across  
70 11 categories spanning four heuristic families (representativeness, cognitive reflection, decision  
71 framing, and loss/availability), with novel isomorphs to resist memorization.
- 72 2. **Behavioral validation** showing that scored lure rates are model-, scorer-, and benchmark-scope  
73 dependent: final-answer rescoring of preserved raw-response artifacts gives lower overall lure  
74 for R1-Distill than Llama (25.5% vs. 33.2%) but high residual vulnerability on base-rate and  
75 conjunction items; OLMo-3-7B shows a stronger Instruct–Think drop (14.9%  $\rightarrow$  0.9%). We  
76 report accuracy and other-response rates so lure suppression is not equated with correctness.
- 77 3. **Probing analysis** showing high conflict/control separability in both models (Llama AUC = 0.974;  
78 R1-Distill AUC = 0.930), with cross-prediction, text baselines, and control tasks used as checks  
79 against purely surface-level explanations.
- 80 4. **Held-in evidence via static probe-direction steering**: in Llama, steering along the probe weight  
81 vector produces a 37.5pp bidirectional behavioral swing; in R1-Distill, the corresponding static  
82 P0 probe direction is linearly decodable (AUC = 0.930) but does not show a stable dose-response  
83 at any of four tested layers (L14/L25/L28/L31) under continuous 2048-token intervention. A  
84 post-hoc P0-only diagnostic on the same 80 conflict items leaves R1 verdicts unchanged at  
85  $\alpha \in \{-5, 0, +5\}$ , but this diagnostic is uncalibrated and near-zero-KL.
- 86 5. **Diagnostic analyses** showing two limits on simple probe narratives: Qwen 3-8B’s hard think/no-  
87 think template yields non-transferring P0 geometries, and within-CoT probe AUC is non-stationary  
88 over a generated trace.

89 Appendix diagnostics also compare available **OLMo 32B** checkpoints, finding higher scored lure  
90 rates for the larger Instruct checkpoint in this family while the Think checkpoint remains near-zero

91 on scored lures. Exploratory SAE and attention entropy analyses are reported as diagnostics rather  
92 than primary evidence.

## 93 2 Related Work

94 **Dual-process cognition in LLMs.** Hagendorff et al. [2023] first documented human-like cognitive  
95 biases in LLMs, finding that instruction-tuned models reproduce representativeness and anchoring  
96 effects that vanish in later ChatGPT versions. Coda-Forno et al. [2025] applied psychometric methods  
97 to show that heuristic and deliberative prompts elicit shared early-layer representations but diverge  
98 in later layers (a suggestive but purely behavioral result). Ziabari et al. [2025] positioned reasoning  
99 along a continuous spectrum rather than a binary, examining how reasoning training shifts the  
100 balance between fast and slow processing modes. Brady et al. [2025] review the dual-process lens  
101 for LLM decision-making in *Nature Reviews Psychology*, noting that LLMs mimic both System 1-  
102 and System 2-like responses but through mechanisms distinct from human cognition. Our work  
103 moves beyond behavioral characterization: we test whether the System 1/System 2 distinction has a  
104 *mechanistic* counterpart in model internals.

105 **Cognitive bias benchmarks and probing.** Malberg et al. [2024] surveyed 30 cognitive biases  
106 across 20 LLMs, establishing that bias susceptibility varies systematically with model scale and  
107 training recipe, but without examining internal representations. Most directly relevant, Huang et al.  
108 [2026] demonstrated that cognitive bias directions are linearly decodable and steerable in LLM  
109 residual streams, achieving 26–32% bias reduction via contrastive activation addition. We differ  
110 from CogBias in three respects: (i) architecturally matched base–reasoning model pairs that hold  
111 architecture fixed while varying the post-training recipe, (ii) Hewitt–Liang control tasks and cross-  
112 domain transfer tests as probe validity checks [Hewitt and Liang, 2019], and (iii) a different causal  
113 question: whether the discriminative probe boundary itself is bidirectionally writable, in addition to  
114 comparing against a CogBias-style mean-difference vector.

115 **Activation steering and representation engineering.** Activation addition [Turner et al., 2023]  
116 and its variants [Panickssery et al., 2024, Li et al., 2024] demonstrate that linear steering vectors can  
117 modulate model behavior at inference time. CogBias [Huang et al., 2026] applied this to cognitive  
118 bias (26–32% reduction) using mean-difference vectors. Probe-direction steering asks whether the  
119 learned discriminative boundary is behaviorally writable; mean-difference activation addition remains  
120 the stronger one-sided mitigation baseline in our Llama control. Pre-CoT probing and steering work  
121 shows that hidden states before generated reasoning can encode answers or downstream reasoning  
122 success [Cox et al., 2026, Afzal et al., 2025, Zhang et al., 2025]; our narrower question is whether  
123 a pre-generation conflict/control direction remains behaviorally writable during reasoning-trained  
124 generation. For SAE analysis, we adopt Ma et al.’s falsification protocol [Ma et al., 2026] as a  
125 mandatory filter, since 45–90% of purported “reasoning features” are spurious token-level artifacts  
126 [Méloux et al., 2025].

127 **Conflict detection in cognitive science.** Our theoretical framework draws on De Neys’ conflict de-  
128 tection paradigm [De Neys, 2012, 2017], which demonstrates that humans detect heuristic–normative  
129 conflicts even when they fail to override the heuristic response. Evans and Stanovich [2013] refined  
130 the Type 1/Type 2 taxonomy, emphasizing that the distinction is graded rather than binary, a framing  
131 we adopt for LLMs. Botvinick et al. [2001]’s conflict monitoring theory provides the computational  
132 account: the anterior cingulate cortex signals response conflict, triggering increased cognitive control.  
133 We test the LLM analogue of this account: do models show internal signatures of conflict detec-  
134 tion (elevated attention entropy, reduced probe confidence) even when producing heuristic-driven  
135 responses?

## 136 3 Benchmark

137 We construct a contrastive cognitive bias benchmark following three design principles: (1) every  
138 conflict item has a structurally matched no-conflict control [De Neys, 2012]; (2) all items use novel  
139 surface features to resist memorization; and (3) category diversity enables specificity analysis.

140 The benchmark comprises **470 items** organized as **235 matched pairs** across 11 categories spanning  
141 four heuristic families: *representativeness* (base rate neglect, 35 pairs; conjunction fallacy, 20; belief-  
142 bias syllogisms, 25), *cognitive reflection* (CRT variants, 30; arithmetic, 25), *decision framing* (framing  
143 effects, 20; anchoring, 20), and *loss/availability* (sunk cost, loss aversion, certainty effect, availability;  
144 60 pairs total). Base rate items include 10 pairs in Gigerenzer’s natural frequency format [Gigerenzer  
145 and Hoffrage, 1995]. Each item is scored as *correct*, *S1-lure* (the specific heuristic-predicted wrong  
146 answer), *other-wrong*, or *refusal*. Classic problems (e.g., the original bat-and-ball) serve only as  
147 contamination baselines.

148 **Item counts and analysis subsets.** Our benchmark contains 470 items across 11 categories.  
149 Primary probe analyses use the 3 vulnerable categories (160 items). Cross-prediction uses all 7  
150 original categories (330 items; 165 pairs). The full 11-category set includes 4 additional categories  
151 added to test domain generality.

## 152 4 Methods

### 153 4.1 Models

154 Our headline comparison exploits the architectural identity between **Llama-3.1-8B-Instruct** [Meta  
155 AI, 2024] and **R1-Distill-Llama-8B** [DeepSeek-AI, 2025]: both have 32 transformer layers, 32 query  
156 heads (8 KV heads via grouped-query attention), and 4096 hidden dimensions. They therefore control  
157 architecture but not the full training intervention: R1-Distill-Llama-8B comes from a different post-  
158 training pipeline, including distillation from a reasoning-trained teacher, CoT formatting, objective  
159 choices, and answer-style changes. We additionally evaluate **Qwen 3-8B** in both thinking (/think)  
160 and non-thinking (/no\_think) modes, providing a same-weights comparison whose hard chat-  
161 template switch changes both reasoning behavior and the P0 positional context. For cross-architecture  
162 replication, we evaluate **OLMo-3-7B-Instruct** and **OLMo-3-7B-Think** [Team OLMo, 2025], which  
163 share the same OLMo base but differ in reasoning training. For scale analysis, we additionally  
164 evaluate **OLMo-3.1-32B-Instruct** and **OLMo-3-32B-Think** (64 layers, 4.5× the parameters of  
165 the 7B variant), enabling a within-OLMo scale comparison on identical benchmark items; because  
166 the 32B Instruct and Think checkpoints use different OLMo release lines, this is not a pure size or  
167 training-objective ablation.

168 **Prompting and decoding.** All models use greedy decoding (`do_sample=False`,  
169 `temperature=1.0`) with no system prompt. For Qwen 3-8B, we use the hard mode switch  
170 (`enable_thinking=True/False` in `apply_chat_template`). In no-think mode, this appends an  
171 empty `<think>\n\n</think>\n\n` block after the assistant turn marker, so the token sequences  
172 *diverge* after that marker: think-mode P0 is the newline after assistant, while no-think-mode P0  
173 follows the injected `</think>` close tag. Cross-mode probe comparisons are therefore template-  
174 conditioned geometry tests, not clean same-position tests of an invariant semantic conflict axis. This  
175 is distinct from Qwen’s soft switch (/think, /no\_think tokens). For R1-Distill, we parse any  
176 generated `<think>...</think>` span, and conservatively treat an open but unclosed thinking span  
177 as no final answer. We use `max_new_tokens=128` for standard models and `max_new_tokens=2048`  
178 for reasoning models to accommodate chain-of-thought traces. We acknowledge this departs from  
179 vendor-recommended settings (e.g., DeepSeek recommends `temperature=0.5-0.7`); Appendix E  
180 reports multi-seed robustness analysis with sampled decoding (`temperature=0.7`) showing  
181 qualitatively consistent results. For R1 steering at the later tested layers (L25/L28/L31), the “other”  
182 response rate is 0% under this parser; the earlier L14 intervention is less fluent, with 30–31.25%  
183 “other” responses.

### 184 4.2 Activation Extraction

185 For each model–item pair, we extract residual stream activations (post-MLP, post-residual) at all  
186 layers at position P0: the last prompt token before generation. This “pre-decision” position captures  
187 the internal state from which the answer is generated. Forward passes use the requested model  
188 precision (bfloat16 where applicable), but HDF5 residual caches store bfloat16 requests as float16  
189 because h5py has no native bfloat16 dtype. For within-CoT probing (Section 5.4), we additionally  
190 extract activations at 10 evenly-spaced token positions within the reasoning trace.

### 191 4.3 Linear Probes with Hewitt–Liang Controls

192 We train  $\ell_2$ -regularized logistic regression probes (LogisticRegressionCV, 5-fold stratified cross-  
193 validation) on residual stream activations at each layer. The probe target is binary: *conflict* vs.  
194 *no-conflict* condition, restricted to the three vulnerable categories (base rate neglect, conjunction  
195 fallacy, syllogisms) where the standard model shows non-trivial lure rates.

196 **Controls.** We apply Hewitt–Liang control tasks [Hewitt and Liang, 2019]: probes trained on  
197 randomly permuted labels establish a selectivity floor. Selectivity (true AUC minus control AUC)  
198 below 5 percentage points indicates the signal reflects probe capacity, not representational structure.  
199 We report ROC-AUC as the primary metric with bootstrap 95% confidence intervals (1000 resamples).

200 **Lower-vulnerability categories.** We separately probe four categories that were lower-vulnerability  
201 in the original source runs (CRT, arithmetic, framing, anchoring) as a *negative control*: if probes  
202 achieve high AUC on categories where there is minimal behavioral difference, the signal is confounded  
203 by surface features rather than processing mode.

### 204 4.4 Probe-Direction Steering

205 To test behavioral control by a static linear probe direction, we extract the weight vector  $\mathbf{w}$  from the  
206 trained logistic regression probe at the peak layer and use it as a steering direction. At inference  
207 time, we add  $\alpha \cdot \hat{\mathbf{w}}$  (where  $\hat{\mathbf{w}} = \mathbf{w}/\|\mathbf{w}\|$ ) to the residual stream at the intervention layer, sweeping  
208 positive values toward S2-like and negative values toward S1-like. The coefficient  $\alpha$  is a unit-direction  
209 residual multiplier, not a perturbation calibrated to equal KL shift or activation-norm change across  
210 models or layers; matched  $\alpha$  sweeps should therefore be read as a fixed-protocol comparison rather  
211 than a matched-strength comparison. We evaluate the steered model on the 80 vulnerable-category  
212 conflict items and report lure-rate changes; specificity controls on matched non-conflict items appear  
213 in Appendix H.

## 214 5 Results

### 215 5.1 Behavioral Validation

216 Table 1 reports lure rates (proportion of S1-lure responses on conflict items). For Llama/R1/Qwen,  
217 we use the preserved 470-item raw-response artifacts and rescore only the final-answer span after  
218 any generated thinking trace; this avoids counting lure mentions inside reasoning text. OLMo rows  
219 use the original full-run verdict artifacts because raw responses were not cached. The compact table  
220 shows selected categories; full per-category rates appear in Appendix G.

221 Four findings emerge. First, behavioral vulnerability is not a simple reasoning/non-reasoning binary.  
222 Under final-answer rescoring of preserved 470-item raw responses, R1-Distill reduces overall lure  
223 relative to Llama (25.5% vs. 33.2%) but remains highly vulnerable on base rate (62.9%) and con-  
224 junction items (80.0%). This replaces the older compact seven-category source-run estimate, whose  
225 raw responses were not fully available for rescoring. Second, vulnerability is *category-specific*: base  
226 rate neglect and conjunction fallacy are consistently susceptible, whereas several other categories are  
227 model- and scorer-sensitive rather than uniformly immune. Third, different models show distinct  
228 profiles (for example, Qwen no-think has 95.0% conjunction lure but 0.0% syllogism lure), suggesting  
229 model-specific rather than universal heuristic behavior. Fourth, Qwen’s within-model thinking toggle  
230 reduces overall lure from 31.1% to 7.7% under strict final-answer rescoring, but mostly by shifting  
231 unfinished traces into the ‘other’ category (56.6%) rather than by pure accuracy gain; conjunction  
232 remains vulnerable at 30.0%. Lure suppression is not automatically normative reasoning: “other”  
233 answers are counted as errors and reported in Appendix G. We therefore interpret lure, accuracy, and  
234 other-response rates jointly rather than treating lure-rate changes as sufficient evidence of reasoning  
235 success.

236 **Cross-architecture replication.** OLMo-3-7B-Instruct/Think [Team OLMo, 2025] independently  
237 show a stronger Instruct–Think drop: 14.9%  $\rightarrow$  0.9% lure rate on the full 470-item run. Probing shows  
238 the same mechanistic direction: OLMo Instruct AUC = 0.996 at L24, OLMo Think AUC = 0.962 at  
239 L22 (non-overlapping CIs), with the standard model showing higher separability in 30/32 layers.

Table 1: Behavioral results on conflict items. Bolded entries indicate vulnerable categories (>10% lure). Llama/R1/Qwen entries are final-answer rescoring of preserved 470-item raw-response artifacts; OLMo entries are the original 470-item verdict artifacts. “Overall” is a micro-average over all 235 conflict items. Behavioral rates are used as validation and sensitivity evidence, not as the causal steering result.

Category	Llama	R1-Distill	Qwen 3-8B		OLMo	OLMo
	8B-Inst.	Llama-8B	no think	think	3-7B-Inst.	3-7B-Think
Base rate	<b>88.6</b>	<b>62.9</b>	<b>71.4</b>	2.9	<b>46</b>	0
Conjunction	<b>55.0</b>	<b>80.0</b>	<b>95.0</b>	<b>30.0</b>	<b>50</b>	0
Syllogism	<b>32.0</b>	<b>20.0</b>	0.0	0.0	0	4
CRT variants	10.0	<b>16.7</b>	<b>13.3</b>	0.0	3	3
Arithmetic	0.0	0.0	0.0	0.0	0	0
Framing	<b>45.0</b>	<b>35.0</b>	<b>45.0</b>	5.0	0	0
Anchoring	10.0	5.0	<b>15.0</b>	0.0	5	0
<b>Overall lure</b>	<b>33.2</b>	<b>25.5</b>	<b>31.1</b>	7.7	<b>14.9</b>	<b>0.9</b>

	Llama	R1	Qwen-NT	Qwen-T	OLMo-I	OLMo-T
Overall Corr	59.6	68.1	59.1	35.7	75	88
Overall Lure	33.2	25.5	31.1	7.7	14.9	0.9
Overall Other	7.2	6.4	9.8	56.6	10	11

240 **Implicit conflict detection.** Llama’s first-token probability is 4.2pp lower on conflict items (0.751  
241 vs. 0.793), and entropy is 16% higher (0.948 vs. 0.817), an implicit “hesitation” signature that  
242 parallels De Neys’ conflict-detection findings [De Neys, 2012] and elevated Anterior Cingulate  
243 Cortex (ACC) activation in humans facing dual-process conflicts [Botvinick et al., 2001]. This is a  
244 pre-generation hesitation signature consistent with conflict sensitivity, even though Llama often fails  
245 to resolve the conflict in its final answer.

## 246 5.2 Linear Probes and Cross-Prediction

247 Figure 1 shows layer-wise probe AUC for the conflict/control classification on the three vulnerable  
248 categories.

249 Llama peaks at AUC = 0.974 (95% CI [0.952, 0.992]) at L16; R1-Distill peaks lower at AUC = 0.930  
250 [0.894, 0.960] at L31 (CIs overlap slightly but point estimates differ by 0.044). The 15-layer  
251 shift (L16 → L31 of 32) is architecturally suggestive but not a localization proof: it shows that the  
252 conflict/control distinction is most linearly accessible much later in R1-Distill, near the hidden state  
253 that seeds the chain-of-thought rollout [Zhang et al., 2025]. Qwen 3-8B achieves high separability in  
254 both modes (no-think: 0.954; think: 0.970; both at L34), but cross-mode probe transfer is at chance  
255 (AUC = 0.496, cosine = -0.005): probes trained in one mode do not generalize to the other. Because  
256 the no-think template injects four tokens before P0 (Section 4), the two modes’ P0 positions differ  
257 in both absolute index and local context, inducing a geometric subspace change through RoPE and  
258 attention to the injected tokens. The non-transferability is therefore an expected consequence of the  
259 template difference, not evidence by itself that the semantic conflict representation is mode-specific.  
260 The useful observation is narrower: each positional context independently supports high within-mode  
261 separability (~0.97), while the hard mode switch prevents a single P0 linear probe from transferring  
262 across templates.

263 **Selectivity and text-only baseline.** Legacy Hewitt–Liang control probes achieve  $AUC \leq 0.66$   
264 at all layers, indicating that the true-probe signal is not explained by probe capacity alone under  
265 that control protocol. Because the repository now uses stricter shuffled-label control semantics,  
266 these selectivity values should be treated as a diagnostic rather than a fresh confirmatory statistic.<sup>1</sup>

<sup>1</sup>We train probes in two settings: 5-fold CV on all 330 items (7 categories, peak L14, AUC = 0.999), used for cross-prediction and steering; and bootstrap CI analysis on 160 vulnerable-category items (3 categories, peak L16, AUC = 0.974 [0.952, 0.992]). The steering direction uses the legacy plain-label-stratified L14 CV probe; the L16 value is the post-selected bootstrap peak used for reporting decodability.

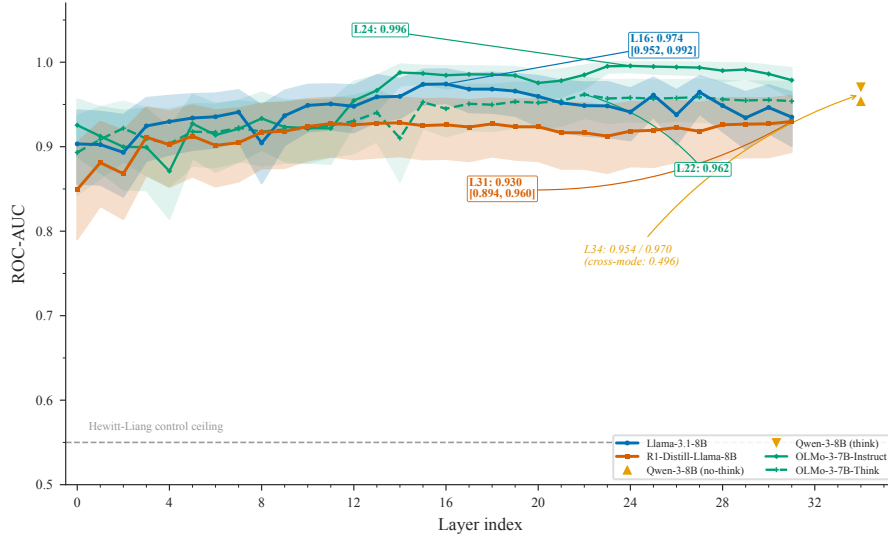


Figure 1: Layer-wise probe ROC-AUC (conflict vs. no-conflict, three vulnerable categories). Llama peaks at 0.974 [0.952, 0.992] (L16); R1-Distill at 0.930 [0.894, 0.960] (L31). Qwen 3-8B peaks at 0.954 (no-think) and 0.970 (think) at L34, but hard-template/position-shifted cross-mode transfer is at chance (0.496). Gray dashed: Hewitt–Liang control ceiling.

Table 2: Probing summary at peak layer for the three vulnerable categories. 95% CIs: bootstrap (1000 resamples) for Llama/R1/OLMo-7B rows; CV-fold  $\pm 1.96\sigma$  for OLMo-32B rows ( $\dagger$ ; per-sample predictions not cached for bootstrap). Qwen CIs omitted (only layer-level aggregates stored). Layer-selection uncertainty is not propagated in any row. Selectivity = true AUC – control AUC. Cross-mode column: AUC when a probe trained on one Qwen mode is tested on the other. Text baselines: three methods on input text only (no activations)—see range in last row.

Model	Peak L	Peak AUC	95% CI	Ctrl AUC	Select.	Cross
Llama-3.1-8B-Instruct	L16	0.974	[0.952, 0.992]	0.63	0.343	—
R1-Distill-Llama-8B	L31	0.930	[0.894, 0.960]	0.55	0.377	—
Qwen 3-8B (no think)	L34	0.954	—	—	—	0.496
Qwen 3-8B (think)	L34	0.970	—	—	—	0.496
OLMo-3-7B-Instruct	L24	0.996	[0.988, 1.000]	—	—	—
OLMo-3-7B-Think	L22	0.962	[0.934, 0.982]	—	—	—
OLMo-3.1-32B-Instruct	L20	0.9999	[0.9997, 1.000] <sup>†</sup>	—	—	—
OLMo-3-32B-Think	L20	0.9978	[0.9960, 0.9996] <sup>†</sup>	—	—	—
<i>Text baselines</i>	—	0.820–0.863	—	—	—	—

267 Three text-only baselines cluster in a narrow band: a length-only probe (word count, 5-fold CV)  
 268 achieves AUC = 0.863, a TF-IDF unigram+bigram logistic regression achieves 0.820, and a global  
 269 sentence-transformer baseline achieves 0.840, all substantially below the combined vulnerable-  
 270 category activation probe (0.999). All three text-only syllogism baselines collapse to below chance  
 271 ( $\leq 0.086$ ), because conflict and control syllogisms are surface-matched in length ( $d = +0.02$ ) and  
 272 vocabulary. By contrast, a syllogism-only activation probe still reaches 0.936, while the combined  
 273 vulnerable-category probe reaches 0.999. The high 0.82–0.86 text-baseline band is itself an important  
 274 caveat: many conflict/control labels are recoverable from prompt surface statistics. Our claim is only  
 275 that the tested residual-stream representation makes the distinction more linearly accessible than these  
 276 text encoders. Since activations are deterministic functions of the input text and model weights, this  
 277 is not evidence of information absent from the prompt, nor does it rule out richer surface-template  
 278 confounds.

279 **Probe summary.** Table 2 summarizes the key probing metrics.

280 **Cross-prediction and immune-category specificity.** Training on *vulnerable* categories and testing  
281 on *immune* categories yields non-transfer (L14 AUC = 0.378). The below-chance AUC reflects a  
282 small polarity reversal inside a compressed immune cluster: immune conflict/control separation is  
283 only  $0.03\times$  the vulnerable separation, while the immune score spread is about  $0.25\times$  as large. Base  
284 rate and conjunction probes transfer near-perfectly to each other (AUC > 0.99), while transfer to  
285 syllogism is weaker (0.59–0.63).

286 **Cross-model probe transfer.** A Llama-trained probe achieves AUC = 0.920 on R1-Distill (layer 23,  
287 the peak transfer layer from a full 32-layer sweep); the reverse achieves 0.954 at layer 15 (also the  
288 peak transfer layer from a full sweep). The processing-mode direction is substantially shared across  
289 this same-architecture pair despite different post-training pipelines, contrasting with Huang et al.  
290 [2026]’s near-orthogonal finding across architectures.

### 291 5.3 Held-in Steering Evidence

292 We steer with the normalized coefficient vector from the vulnerable-category L14 logistic-regression  
293 probe (CV AUC = 0.960), adding  $\alpha\bar{w}$  to the residual stream for  $\alpha \in \{-5, -3, -1, 0, +1, +3, +5\}$  on  
294 the 80 vulnerable conflict items. This is a held-in behavioral test: the probe is fit on the same 160-item  
295 family, so the causal claim rests on direction-specific controls rather than held-out generalization.

296 Figure 2 shows a monotonic Llama response: lure rate moves from 68.8% at  $\alpha = -5$  to 31.2% at  
297  $\alpha = +5$  (37.5pp). The entire lure reduction converts to correct answers (other rate stays 0), and  
298 50 random unit directions show no systematic endpoint effect (mean lure 52–57%,  $\sigma \leq 8.7$ pp).  
299 Two controls sharpen the interpretation. Probe steering produces zero lure responses on matched  
300 non-conflict items, although large steering can still degrade control accuracy (Appendix H). The  
301 CogBias-style mean-difference vector is a stronger one-sided debiaser (52.5%  $\rightarrow$  11.25% lure at  
302  $\alpha = +5$ ) but is not bidirectional; negative  $\alpha$  leaves lure at 53.75%. The probe direction is weaker for  
303 debiasing but better evidence of a behaviorally effective boundary under this protocol.

304 **R1-Distill steering contrast: layer sweep.** We sweep R1-Distill layers L14/L25/L28/L31 under  
305 continuous 2048-token steering in the original source-run scoring regime (Appendix A). No tested  
306 layer shows a coherent dose-response: lure fluctuates between 5–14%, and at the probe peak L31 it  
307 changes from 10.0% to 11.25% between  $\alpha = -5$  and  $+5$  (5.0pp full range), versus Llama’s 37.5pp  
308 endpoint swing. The direction that controls Llama therefore does not yield a stable dose-response in  
309 R1 under this static continuous-addition protocol, and no matched R1 random-direction control band  
310 is included for this source-run sweep. However, the corrected final-answer-span audit changes the R1  
311 baseline and is the stricter diagnostic.

312 As an exploratory post-hoc position-targeted check, we add the L31 vector only at P0 and disable  
313 steering during generation. Across all 80 vulnerable conflict items, R1 gives identical verdicts  
314 at  $\alpha = -5, 0, +5$  (56.2% lure; paired endpoint CI [0.0, 0.0]pp), and the pattern also holds at  
315  $\alpha = \pm 20$ . This check argues against the narrow hypothesis that the continuous sweep failed only  
316 because later-token steering washed out the P0 perturbation. It is not a clean held-out category-  
317 transfer claim because a probe-check showed below-chance held-out projection transfer for syllogism;  
318 moreover,  $\alpha = \pm 5$  prompt-prefill KL is  $< 10^{-6}$  and matched all-vulnerable R1 random-direction  
319 controls are not included. A narrower base-rate/conjunction P0-only random-control audit gives the  
320 same 0pp endpoint swing for the probe and 10 random directions, but it shares the near-zero-KL  
321 limitation (Appendix A). A completed high-alpha base-rate/conjunction diagnostic instead produces  
322 macroscopic prompt-logit KL and changes R1 behavior, reinforcing that the main result is a unit-scale  
323 transfer failure rather than evidence of general non-writability. The comparison is therefore not a  
324 matched-perturbation test: Llama shows a large behavioral swing, whereas the R1 P0 diagnostic has  
325 near-zero prompt-logit impact. Thus the supported claim is a static unit-scale steering-transfer failure:  
326 linear readability persists, but this specific protocol does not provide matched behavioral control. It is  
327 not a claim against dynamic, calibrated multi-layer, SAE-feature, or learned nonlinear interventions.

### 328 5.4 Within-CoT Probing

329 At Qwen 3-8B L34, seven-position probes trace a non-monotonic path (Figure 3): P0 0.938  $\rightarrow$  T0  
330 0.973  $\rightarrow$  T75 0.754  $\rightarrow$  Tend 0.971; P2 returns to chance (0.500). We interpret this conservatively:  
331 the boundary is non-stationary over the trace, but the AUC changes could reflect boundary rotation,

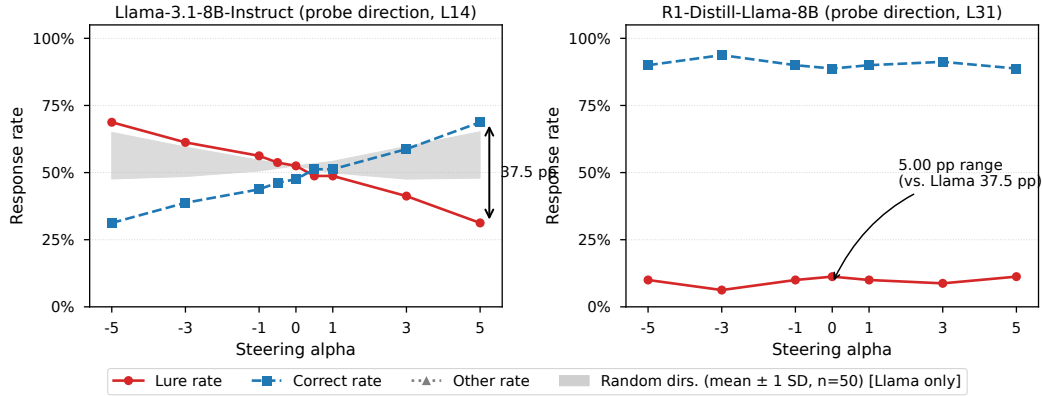


Figure 2: Probe-direction steering (80 vulnerable-category conflict items). *Left*: Llama-3.1-8B-Instruct (L14) shows a 37.5pp dose-response ( $\alpha = -5$ : 68.8% lure;  $\alpha = +5$ : 31.2%), with 100% of lure reduction converting to correct answers (other rate stays 0); 50 random directions (gray band,  $n = 50$ ) show no systematic effect. *Right*: R1-Distill-Llama-8B (L31, its probe-peak layer) shows no coherent dose-response (5.0pp range, compared to Llama’s 37.5pp) under the original static continuous-addition source-run scoring protocol.

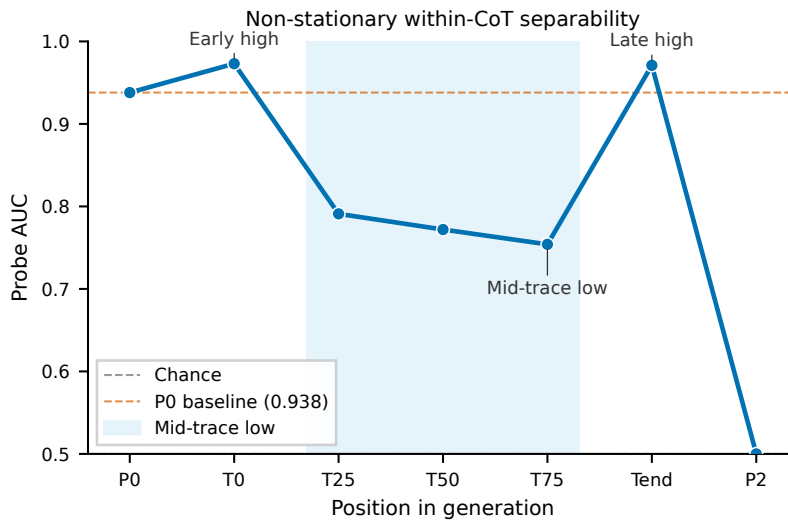


Figure 3: Within-CoT probe AUC for Qwen 3-8B (think mode) at L34. The trajectory shows an early rise (P0→T0, +3.5pp), a mid-trace dip (T0→T75, -22pp), and a late rebound (T75→Tend, 0.971). P2 (control) = 0.500.

332 dispersion, position effects, or other context-dependent shifts. They do not establish stepwise  
 333 reasoning or causal use of specific CoT tokens [Boppana et al., 2026, Zhao et al., 2025]. Full  
 334 trajectories in Appendix B.

## 335 References

- 336 Anum Afzal, Florian Matthes, Gal Chechik, and Yftah Ziser. Knowing before saying: LLM  
337 representations encode information about chain-of-thought success before completion. In *Findings*  
338 *of the Association for Computational Linguistics: ACL 2025*, 2025.
- 339 Siddharth Boppana, Annabel Ma, Max Loeffler, Raphael Sarfati, Eric Bigelow, Atticus Geiger, Owen  
340 Lewis, and Jack Merullo. Reasoning theater: Disentangling model beliefs from chain-of-thought.  
341 *arXiv preprint arXiv:2603.05488*, 2026.
- 342 Matthew M Botvinick, Todd S Braver, Deanna M Barch, Cameron S Carter, and Jonathan D Cohen.  
343 Conflict monitoring and cognitive control. *Psychological Review*, 108(3):624–652, 2001.
- 344 Oliver Brady, Paul Nulty, Lili Zhang, Tomas Ward, and David McGovern. Dual-process theory and  
345 decision-making in large language models. *Nature Reviews Psychology*, 4:777–792, 2025. doi:  
346 10.1038/s44159-025-00506-1.
- 347 Julian Coda-Forno, Zhuokai Zhao, Qiang Zhang, Dipesh Tamboli, Weiwei Li, Xiangjun Fan, Lizhu  
348 Zhang, Eric Schulz, and Hsiao-Ping Tseng. Exploring system 1 and 2 communication for latent  
349 reasoning in LLMs. *arXiv preprint arXiv:2510.00494*, 2025.
- 350 Kyle Cox, Darius Kianersi, and Adrià Garriga-Alonso. Decoding answers before chain-of-thought:  
351 Evidence from pre-CoT probes and activation steering. *arXiv preprint arXiv:2603.01437*, 2026.
- 352 Wim De Neys. Bias and conflict: A case for logical intuitions. *Perspectives on Psychological Science*,  
353 7(1):28–38, 2012.
- 354 Wim De Neys, editor. *Dual Process Theory 2.0*. Routledge, London, 2017.
- 355 Wim De Neys. Advancing theorizing about fast-and-slow thinking. *Behavioral and Brain Sciences*,  
356 46:e111, 2023.
- 357 DeepSeek-AI. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning.  
358 *arXiv preprint arXiv:2501.12948*, 2025.
- 359 Jonathan St BT Evans and Keith E Stanovich. Dual-process theories of higher cognition: Advancing  
360 the debate. *Perspectives on Psychological Science*, 8(3):223–241, 2013.
- 361 Gerd Gigerenzer and Ulrich Hoffrage. How to improve Bayesian reasoning without instruction:  
362 Frequency formats. *Psychological Review*, 102(4):684–704, 1995.
- 363 Thilo Hagendorff, Sarah Fabi, and Michal Kosinski. Human-like intuitive behavior and reasoning  
364 biases emerged in large language models but disappeared in ChatGPT. *Nature Computational*  
365 *Science*, 3(10):833–838, 2023.
- 366 John Hewitt and Percy Liang. Designing and interpreting probes with control tasks. In *Proceedings*  
367 *of the 2019 Conference on Empirical Methods in Natural Language Processing*, pages 2733–2743,  
368 2019.
- 369 Fan Huang, Songheng Zhang, Haewoon Kwak, and Jisun An. CogBias: Measuring and mitigating  
370 cognitive bias in large language models. *arXiv preprint arXiv:2604.01366*, 2026.
- 371 Daniel Kahneman. *Thinking, Fast and Slow*. Farrar, Straus and Giroux, New York, 2011.
- 372 Su Hwan Kim, Sebastian Ziegelmayr, Felix Busch, Christian J Mertens, Matthias Keicher, Lisa C  
373 Adams, Keno K Bressen, Rickmer Braren, Marcus R Makowski, Jan S Kirschke, Dennis M  
374 Hedderich, and Benedikt Wiestler. LLM reasoning does not protect against clinical cognitive  
375 biases—an evaluation using BiasMedQA. *medRxiv*, 2025.
- 376 Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Her-  
377 nandez, Dustin Li, Esin Durmus, Evan Hubinger, Deep Ganguli, et al. Measuring faithfulness in  
378 chain-of-thought reasoning. *arXiv preprint arXiv:2307.13702*, 2023.
- 379 Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time  
380 intervention: Eliciting truthful answers from a language model. *Advances in Neural Information*  
381 *Processing Systems*, 36, 2024.

- 382 George Ma, Zhongyuan Liang, Irene Y. Chen, and Somayah Sojoudi. Do sparse autoencoders identify  
383 reasoning features in language models? *arXiv preprint arXiv:2601.05679*, 2026.
- 384 Simon Malberg, Roman Poletukhin, Carolin M Schuster, and Georg Groh. A comprehensive  
385 evaluation of cognitive biases in LLMs. *arXiv preprint arXiv:2410.15413*, 2024.
- 386 David E Melnikoff and John A Bargh. The mythical number two. *Trends in Cognitive Sciences*, 22  
387 (4):280–293, 2018.
- 388 Maxime M eloux, Giada Dirupo, Franois Portet, and Maxime Peyrard. The dead salmon of AI  
389 interpretability. *arXiv preprint arXiv:2512.18792*, 2025.
- 390 Meta AI. The Llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 391 Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt  
392 Turner. Steering Llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*,  
393 2024.
- 394 Team OLMo. OLMo 3. *arXiv preprint arXiv:2512.13961*, 2025.
- 395 Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J Vazquez, Ulisse Mini,  
396 and Monte MacDiarmid. Steering language models with activation engineering. *arXiv preprint*  
397 *arXiv:2308.10248*, 2023.
- 398 Anqi Zhang, Yulin Chen, Jane Pan, Chen Zhao, Aurojit Panda, Jinyang Li, and He He. Reasoning  
399 models know when they’re right: Probing hidden states for self-verification. *arXiv preprint*  
400 *arXiv:2504.05419*, 2025.
- 401 Jiachen Zhao, Yiyu Sun, Weiyan Shi, and Dawn Song. Can aha moments be fake? identifying true  
402 and decorative thinking steps in chain-of-thought. *arXiv preprint arXiv:2510.24941*, 2025.
- 403 Alireza S Ziabari, Nona Ghazizadeh, Zhivar Sourati, Farzan Karimi-Malekabadi, Payam Piray, and  
404 Morteza Dehghani. Reasoning on a spectrum: Aligning LLMs to system 1 and system 2 thinking.  
405 *arXiv preprint arXiv:2502.12470*, 2025.

## 406 A R1 Continuous Layer-Sweep Steering Details

Table 3: R1-Distill layer-sweep steering (continuous, 80 vulnerable conflict items, original source-run scoring). Lure rate (%) at each  $\alpha$ ; none of the four tested layers shows a monotonic dose-response (full-sweep ranges: L14 = 8.8pp, L25 = 7.6pp, L28 = 5.0pp, L31 = 5.0pp; endpoint swing at L31: 10.0%  $\rightarrow$  11.25% = 1.25pp). Compare Llama L14: 68.8%  $\rightarrow$  31.2% (37.5pp endpoint swing).

Layer	-5	-3	-1	0	+1	+3	+5
L14	6.2	10.0	7.5	6.2	5.0	7.5	13.8
L25	6.2	6.2	8.8	11.2	10.0	13.8	8.8
L28	8.8	6.2	10.0	11.2	7.5	11.2	10.0
L31	10.0	6.2	10.0	11.2	10.0	8.8	11.2

407 **Position-targeted diagnostics.** Table 4 reports exploratory post-hoc R1 diagnostics that perturb  
408 fewer positions than the continuous layer sweep. These runs are useful checks against the simplest  
409 washout objection, but they are not confirmatory: matched R1 random-direction controls for the  
410 continuous and position-targeted runs are not included, the P0 KL shift is tiny at  $|\alpha| = 5$ , the all-  
411 vulnerable rows include syllogism despite below-chance held-out projection transfer in a probe-check,  
412 and answer-after-think covers only the 55 base-rate/conjunction conflict items.

413 A narrower unit-scale random-control audit on the 55 base-rate/conjunction conflict items gives  
414 69.1% lure at  $\alpha = -5, 0, +5$  for the probe direction, and 10 matched random directions also have  
415 0pp endpoint swing (2.5–97.5% range [0.0, 0.0]pp). Because prompt-prefill KL remains  $\leq 2 \times 10^{-6}$   
416 and syllogism is excluded, this is a random-direction sanity check for the unit-scale diagnostic, not  
417 evidence of calibrated R1 resistance.

Table 4: R1-Distill position-targeted diagnostics. Entries are conflict-item lure rates (%); all rows use the L31 probe direction except the explicit L25+L28+L31 multilayer row. The final column reports mean prompt-prefill KL at  $\alpha = +5$  when that position is perturbed; n/a means the protocol does not perturb prompt-prefill logits.

Protocol	Items	$\alpha = -5$	$\alpha = 0$	$\alpha = +5$	KL@+5
P0 only, L31	80	56.2	56.2	56.2	$4.0 \times 10^{-7}$
T0 only, L31	80	56.2	56.2	56.2	n/a
P0 only, L25+L28+L31	80	58.8	56.2	50.0	$5.3 \times 10^{-6}$
Answer-after-think, L31	55	70.9	69.1	65.5	n/a

418 **High-alpha P0 diagnostic.** Table 5 reports a completed, post-hoc base-rate/conjunction P0-only  
 419 diagnostic that deliberately leaves the unit-scale regime. The result is not a calibrated confirmation of  
 420 a mechanism: it has no random-direction controls, covers 55 conflict/control pairs rather than all three  
 421 vulnerable categories, and uses very large perturbations whose prompt-prefill KL is around 10–13  
 422 with a 100% top-token flip rate. It is nevertheless important for scope. At  $\alpha = -500$  or  $-1000$ , R1  
 423 lure falls from 69.1% to 18.2%; at  $\alpha = +500$  or  $+1000$ , lure falls to 52.7%. Thus R1 is not claimed  
 424 to be unwritable in general. The paper’s main claim is that the unit-scale static direction that controls  
 425 Llama does not transfer as a matched behavioral-control protocol to R1-Distill.

Table 5: R1-Distill high-alpha P0 diagnostic on base-rate/conjunction conflict items ( $n = 55$ ), final-answer-span scoring. These very large perturbations produce macroscopic prompt-prefill KL and change behavior, delimiting the unit-scale claim rather than supporting a general non-writability conclusion.

$\alpha$	lure	correct	other	mean KL
-1000	18.2	81.8	0.0	11.49
-500	18.2	81.8	0.0	10.09
0	69.1	30.9	0.0	0.00
+500	52.7	45.5	1.8	11.88
+1000	52.7	45.5	1.8	13.30

## 426 B Within-CoT Probing Details

427 This appendix details the within-chain-of-thought probing protocol for Qwen 3-8B (think mode).  
 428 For each model–item pair, we extract residual stream activations at five canonical positions within  
 429 the <think>...</think> trace: T0 (first token), T25, T50, T75, and Tend (final token before  
 430 </think>). We train independent  $\ell_2$ -regularized logistic regression probes (5-fold stratified CV) at  
 431 each (layer, position) combination on the conflict vs. no-conflict classification ( $n = 160$ ).

Table 6: Within-CoT probe AUC trajectory for Qwen 3-8B (think mode) at L34, 7 positions.

Position	Description	Probe AUC
P0	Pre-generation	0.938
T0	Thinking onset	<b>0.973</b>
T25	25% of trace	0.791
T50	Midpoint	0.772
T75	75% of trace	0.754
Tend	Final thinking token	0.971
P2	Control (chance)	0.500

432 The U-shaped trajectory (T0 peak, T25–T75 dip, Tend rebound) is inconsistent with a flat trace and  
 433 with monotonic sharpening, but it is not by itself a causal trace-faithfulness test. The orthogonal P0  
 434 directions between Qwen’s think/no-think modes (cross-mode AUC = 0.496) should be interpreted  
 435 with the hard-template confound in mind: the result shows that template-specific P0 geometries do

436 not transfer linearly, while the within-trace curve shows that the think-mode representation changes  
437 substantially before the answer.

## 438 C Scale Analysis: Full Results

439 OLMo-3.1-32B-Instruct and OLMo-3.1-32B-Think (64 layers, 5120 hidden dims) were evaluated on  
440 the full 470-item benchmark. Lure rate rises from 14.9% (7B) to 19.6% (32B Instruct), with base  
441 rate neglect doubling (46%→74.3%). The 32B Think checkpoint has 0.4% scored lure rate. Probes  
442 achieve AUC = 0.9999 (Instruct L20) and 0.9978 (Think L20): the behavioral gap widens (19.2pp at  
443 32B vs. 14.0pp at 7B) while the representational gap shrinks (0.034 → 0.0021).

## 444 D SAE Features and Attention Entropy

445 We apply a Goodfire SAE (65,536 features) to Llama-3.1-8B-Instruct L19 P0 residual activations.  
446 Reconstruction fidelity is moderate (MSE 0.011, explained variance 74.0%, mean L0 47.9 active  
447 features), so we treat feature findings as exploratory supporting evidence. After BH-FDR correction  
448 across features, 41 features show significant conflict/control activation differences; all 41 pass the Ma  
449 et al. token-injection falsification protocol (0 spurious under our threshold).

450 For attention, we test all 1,024 heads per model using rank-biserial correlation between conflict and  
451 control items, BH-FDR corrected at  $q < 0.05$ . Llama has 473 significant heads and 30 S2-like-  
452 specialized heads; R1-Distill has 515 significant heads and 57 S2-like-specialized heads. At the more  
453 conservative KV-group granularity, the counts are 13/256 for Llama and 24/256 for R1-Distill. These  
454 analyses support the representation-level story but are not causal evidence.

## 455 E Multi-Seed Robustness

456 Multi-seed evaluation with sampled decoding (temperature 0.7) on the 11-category robustness set  
457 shows that Llama’s sampled lure rate is stable across seeds ( $27.5\% \pm 1.3\text{pp}$ ) while R1-Distill’s is  
458 also stable ( $12.1\% \pm 1.0\text{pp}$ ) in the current robustness artifacts, but category profiles shift and these  
459 sampled artifacts have not been regenerated as a fresh final-answer-rescored behavioral table. Probe  
460 results at P0 are unaffected by generation strategy, as they measure the pre-decision representation.

### 461 E.1 Exact Generation Parameters

462 Table 7 reports the exact generation keyword arguments used for each model family.

Table 7: Generation keyword arguments for all models. “Greedy” denotes the main-text configuration;  
“Sampled” denotes the robustness configuration used in this appendix.

Parameter	Greedy (main text)	Sampled (this appendix)
do_sample	False	True
temperature	1.0	0.7
top_p	—	0.95
max_new_tokens	128 (standard) / 2048 (reasoning)	same
System prompt	none	none

463 **Vendor-recommended settings.** Our greedy configuration departs from several vendor recom-  
464 mendations: DeepSeek recommends temperature=0.5-0.7 for R1 models; Qwen recommends  
465 temperature=0.6 with top\_p=0.95 for Qwen 3. We use deterministic greedy decoding throughout  
466 the main text to (a) ensure exact reproducibility without seed dependence, (b) isolate representational  
467 differences from generation-strategy effects, and (c) simplify the probing analysis, since P0 activations  
468 are deterministic given the prompt regardless of downstream sampling. This appendix demonstrates  
469 that the qualitative findings— in particular the S1/S2 lure-rate gap and the stability of probe AUC at  
470 P0—are preserved under vendor-recommended sampled decoding (temperature=0.7).

471 **Robustness argument.** The key robustness claim is twofold: (1) probe results at P0 are *invariant* to  
 472 generation strategy because P0 activations depend only on the prompt, not on downstream sampling;  
 473 and (2) behavioral lure-rate gaps between S1 and S2 models are qualitatively preserved under sampled  
 474 decoding, though the magnitudes shift for reasoning models whose chain-of-thought is sensitive to  
 475 stochastic token selection.

## 476 F Token-Length Confound Analysis

477 Near-perfect probe AUC (e.g., OLMo-32B at 0.9999) can in principle arise from trivial confounds  
 478 like systematic prompt-length differences between conflict and control items. Table 8 reports subword  
 479 token counts per category using the Llama-3.1-8B tokenizer.

Table 8: Prompt subword length (Llama-3.1-8B tokenizer) for conflict vs. control items by category. “*d*” is Cohen’s *d*; positive values indicate conflict items are longer.

Category	Conflict	Control	Cohen’s <i>d</i>	<i>p</i>
Base rate	98.0	66.7	+3.69	$< 10^{-4}$
Conjunction	87.7	63.4	+6.36	$< 10^{-4}$
Syllogism	93.3	93.2	+0.02	0.953
CRT	47.0	51.8	-0.44	0.097
Arithmetic	70.9	51.0	+6.21	$< 10^{-4}$
Framing	97.2	106.2	-4.86	$< 10^{-4}$
Anchoring	60.4	45.5	+4.62	$< 10^{-4}$
<b>Overall</b>	<b>81.1</b>	<b>68.4</b>	+0.45	$< 10^{-4}$

480 **Why length is unlikely to be the primary signal.** Length differences are substantial on base  
 481 rate and conjunction but balanced on syllogism and reversed on framing. Three observations argue  
 482 against length as the dominant confound. First, a probe trained on syllogism alone (where length is  
 483 balanced,  $d = +0.02$ ) achieves within-category AUC = 0.936 and transfers to base rate (0.950) and  
 484 conjunction (0.990)—the shared direction cannot be length. Second, length direction is inconsistent  
 485 across categories, so a pure-length probe would produce anti-correlated cross-category predictions  
 486 rather than the immune-cluster compression we observe. Third, the text baseline, which has direct  
 487 access to length as a feature, caps at 0.840; the 16 pp gap to the activation probe (0.999) is not  
 488 accounted for by these text baselines alone, though richer surface-template features remain a possible  
 489 confound. For OLMo-32B specifically, probe AUC = 0.9999 on the three vulnerable categories  
 490 (including length-balanced syllogism) implies the probe separates syllogism conflict vs. control  
 491 despite near-identical prompt lengths (93.3 vs. 93.2 tokens).

## 492 G Final-Answer Behavioral Rescoring Details

493 Table 9 gives the category-level final-answer rescoring used for the Llama/R1/Qwen rows in Table 1.  
 494 The rescoring recomputes verdicts from preserved 470-item raw responses after removing generated  
 495 thinking text; close-only R1-style continuations containing “`</think>`” without generated “`<think>`”  
 496 are scored only after the close tag. This table replaces the older legacy source-run behavioral table  
 497 for these models. OLMo rows in Table 1 still use their source-run verdict artifacts because the  
 498 corresponding raw-response artifacts were not cached for final-answer rescoring.

Table 9: Final-answer rescoring of preserved 470-item raw-response artifacts. For Llama/R1/Qwen, verdicts are recomputed from the final-answer span after removing generated thinking text. Acc, Lure, and Other are percentages on conflict items; Ctrl is matched-control accuracy.

Model	Category	<i>n</i>	Acc	Lure	Other	Ctrl
Llama	base_rate	35	11.4	88.6	0.0	100.0
Llama	conjunction	20	45.0	55.0	0.0	100.0
Llama	syllogism	25	68.0	32.0	0.0	52.0
Llama	CRT	30	56.7	10.0	33.3	83.3
Llama	arithmetic	25	100.0	0.0	0.0	100.0
Llama	framing	20	55.0	45.0	0.0	100.0

Model	Category	$n$	Acc	Lure	Other	Ctrl
Llama	anchoring	20	55.0	10.0	35.0	75.0
Llama	availability	15	86.7	13.3	0.0	40.0
Llama	certainty_effect	15	26.7	73.3	0.0	100.0
Llama	loss_aversion	15	100.0	0.0	0.0	100.0
Llama	sunk_cost	15	93.3	6.7	0.0	0.0
<b>Llama</b>	<b>Overall</b>	<b>235</b>	<b>59.6</b>	<b>33.2</b>	<b>7.2</b>	<b>80.4</b>
R1-Distill	base_rate	35	37.1	62.9	0.0	100.0
R1-Distill	conjunction	20	20.0	80.0	0.0	100.0
R1-Distill	syllogism	25	80.0	20.0	0.0	80.0
R1-Distill	CRT	30	66.7	16.7	16.7	100.0
R1-Distill	arithmetic	25	96.0	0.0	4.0	100.0
R1-Distill	framing	20	65.0	35.0	0.0	95.0
R1-Distill	anchoring	20	60.0	5.0	35.0	55.0
R1-Distill	availability	15	86.7	13.3	0.0	60.0
R1-Distill	certainty_effect	15	100.0	0.0	0.0	100.0
R1-Distill	loss_aversion	15	73.3	13.3	13.3	93.3
R1-Distill	sunk_cost	15	100.0	0.0	0.0	0.0
<b>R1-Distill</b>	<b>Overall</b>	<b>235</b>	<b>68.1</b>	<b>25.5</b>	<b>6.4</b>	<b>84.3</b>
Qwen think	base_rate	35	37.1	2.9	60.0	82.9
Qwen think	conjunction	20	5.0	30.0	65.0	60.0
Qwen think	syllogism	25	64.0	0.0	36.0	68.0
Qwen think	CRT	30	3.3	0.0	96.7	10.0
Qwen think	arithmetic	25	12.0	0.0	88.0	92.0
Qwen think	framing	20	0.0	5.0	95.0	55.0
Qwen think	anchoring	20	60.0	0.0	40.0	55.0
Qwen think	availability	15	53.3	26.7	20.0	33.3
Qwen think	certainty_effect	15	40.0	40.0	20.0	66.7
Qwen think	loss_aversion	15	100.0	0.0	0.0	53.3
Qwen think	sunk_cost	15	60.0	0.0	40.0	0.0
<b>Qwen think</b>	<b>Overall</b>	<b>235</b>	<b>35.7</b>	<b>7.7</b>	<b>56.6</b>	<b>54.9</b>
Qwen no-think	base_rate	35	28.6	71.4	0.0	100.0
Qwen no-think	conjunction	20	5.0	95.0	0.0	100.0
Qwen no-think	syllogism	25	48.0	0.0	52.0	52.0
Qwen no-think	CRT	30	60.0	13.3	26.7	100.0
Qwen no-think	arithmetic	25	100.0	0.0	0.0	100.0
Qwen no-think	framing	20	55.0	45.0	0.0	100.0
Qwen no-think	anchoring	20	75.0	15.0	10.0	65.0
Qwen no-think	availability	15	100.0	0.0	0.0	60.0
Qwen no-think	certainty_effect	15	13.3	86.7	0.0	100.0
Qwen no-think	loss_aversion	15	100.0	0.0	0.0	100.0
Qwen no-think	sunk_cost	15	100.0	0.0	0.0	0.0
<b>Qwen no-think</b>	<b>Overall</b>	<b>235</b>	<b>59.1</b>	<b>31.1</b>	<b>9.8</b>	<b>83.0</b>

## 499 H Steering Specificity Controls

500 Table 10 reports specificity controls for Llama-3.1-8B-Instruct at L14. Four sweeps: (1) probe-  
501 direction on conflict items (main result); (2) probe-direction on matched non-conflict control items  
502 (null-effect specificity); (3) mean-difference direction on conflict items (CogBias baseline); (4) mean-  
503 difference on control.  $\cos(\text{probe}, \text{mean-diff}) = 0.314$ : the two directions differ substantially.

Table 10: Steering specificity controls on Llama-3.1-8B-Instruct, L14, continuous steering. Conflict items carry an S1 lure answer; control items carry only the normative answer.

Direction	Items	$\alpha = -5$	$-3$	$-1$	$0$	$+1$	$+3$	$+5$
<i>Lure rate (%)</i>								
Probe	Conflict	68.75	61.25	56.25	52.50	48.75	41.25	31.25
Probe	Control	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Mean-diff	Conflict	53.75	55.00	55.00	52.50	36.25	21.25	11.25
Mean-diff	Control	0.00	0.00	0.00	0.00	0.00	0.00	0.00
<i>Correct rate (%)</i>								
Probe	Conflict	31.25	38.75	43.75	47.50	51.25	58.75	68.75
Probe	Control	78.75	87.50	98.75	100.00	100.00	100.00	100.00
Mean-diff	Conflict	46.25	45.00	45.00	47.50	63.75	78.75	88.75
Mean-diff	Control	60.00	73.75	97.50	100.00	100.00	100.00	100.00

504 **Three observations.** (i) The probe direction produces *zero* lure responses on control items at any  
505  $\alpha$ , so it does not induce the specific heuristic-lure answer on matched no-conflict prompts. The small  
506 drop in correctness at  $\alpha = -5$  on control items (78.75%) reflects residual-stream coherence loss (the  
507 21.25% of items become “other” responses, not lures), the same pattern seen with random directions  
508 in the main text. (ii) The probe-weight direction is *bidirectional*: moving toward S1 (negative  $\alpha$ )  
509 increases lure by 16.25pp, moving toward S2 (positive  $\alpha$ ) decreases lure by 21.25pp. Bidirectional  
510 dose-response is stronger evidence that this direction controls the measured lure behavior than a  
511 unidirectional effect. (iii) The mean-difference baseline is *unidirectional*: negative  $\alpha$  has essentially  
512 no effect (lure stays around 53–55%), while positive  $\alpha$  strongly reduces lure (to 11.25% at  $\alpha = +5$ ).  
513 Mean-difference steering efficiently improves behavior but cannot test causal sufficiency by worsening  
514 it—a key methodological motivation for probe-weight steering.

## 515 I Limitations and Safety Note

516 The main claim is restricted to the tested static continuous and limited position-targeted unit-scale  
517 probe-direction protocols. The completed high-alpha P0 diagnostic shows that much larger perturba-  
518 tions can change R1 behavior, so the result should not be read as non-writability in general. The R1  
519 continuous sweep and position-targeted diagnostics still lack complete matched random-direction  
520 controls, strong cross-model perturbation calibration, sampled-decoding steering, and dynamic or  
521 SAE-feature interventions. Linear probes may miss nonlinear structure; Qwen cross-mode transfer  
522 is confounded by the hard template’s positional shift; within-CoT phase labels are descriptive; and  
523 the benchmark covers English-language cognitive-bias items on models up to 32B parameters. The  
524 safety implication is cautionary: steering safeguards validated on instruction-tuned models should be  
525 revalidated before deployment on reasoning-distilled models. We release no new model weights or  
526 personal data.