
Towards a Mechanistic Understanding of Robustness in Finetuned Reasoning Models

Abstract

1 Supervised fine-tuning (SFT) on chain-of-thought data induces brittleness in lan-
2 guage models, improving reasoning capabilities while severely degrading general
3 performance. We provide the first mechanistic explanation for this trade-off through
4 three complementary techniques: crosscoders for mapping feature transformations,
5 Fisher Information-based identification of causal features, and gradient blocking
6 for intervention experiments. Our analysis reveals that SFT operates through
7 two distinct mechanisms—repurposing shared features for reasoning tasks and
8 suppressing base-only features. Fisher Information with Sparse Autoencoders
9 identifies the specific features responsible for reasoning, validated through fea-
10 ture steering that achieves 3.46% performance gains on base models. Crosscoder
11 analysis demonstrates that SFT repurposes existing reasoning capabilities in the
12 base model rather than creating new ones. Gradient blocking experiments prove
13 these mechanisms are separable: blocking shared features eliminates reasoning
14 entirely, while blocking base-only features preserves it, demonstrating that base
15 feature suppression is unnecessary for reasoning. This mechanistic understanding
16 provides the foundation for developing surgical training methods that preserve
17 general capabilities while enhancing reasoning.

18 1 Introduction

19 Supervised fine-tuning (SFT) on chain-of-thought data is the standard method for enhancing reasoning
20 in language models [8], yet it systematically induces brittleness: mathematical reasoning in Qwen3-
21 4B improves by 62%, but non-reasoning capabilities collapse by 47%. This trade-off persists across
22 model families, leaving practitioners reliant on solutions like KL regularization or RLVR whose
23 mechanisms remain opaque. We provide the first mechanistic explanation for SFT-induced brittleness.
24 Using crosscoders [13], Fisher Information-based feature identification, and a novel intervention
25 called gradient blocking, we discover that SFT operates through two distinct mechanisms: (1)
26 *learning* through repurposing existing shared features, and (2) *suppression* of base-only features—an
27 unnecessary side effect. Our contributions include:

- 28 1. **Mechanistic characterization of SFT.** Using crosscoders, we discover that SFT operates through
29 two distinct processes: repurposing shared features for reasoning, and suppressing base-only
30 features. This provides the first mechanistic explanation for SFT-induced brittleness.
- 31 2. **Identification and localization of reasoning features.** We develop a Fisher Information-based
32 method with Sparse Autoencoders that identifies reasoning features. Feature steering validates
33 these are causal, achieving 3.46% performance gains on the base models and outperforming exist-
34 ing feature identification methods. Crosscoder analysis reveals these features maintain their direc-
35 tion after finetuning, proving that SFT repurposes existing features rather than creating new ones.
- 36 3. **Causal proof of mechanism separability via gradient blocking.** We introduce gradient
37 blocking to selectively freeze feature subsets during training. Blocking shared features eliminates
38 reasoning entirely, proving their modification is necessary for learning. Blocking base-only
39 features preserves reasoning, suggesting their suppression is unnecessary for reasoning.

40 2 Methodology

41 **Identifying Causal Reasoning Features.** A mechanistic account of how SFT affects reasoning
42 requires identification of the specific, interpretable features that constitute this capability. We employ
43 SAEs [2] to decompose model activations $h \in \mathbb{R}^d$ into sparse feature representations, where $f_j(h)$
44 denotes the activation of feature j . Following theoretical insights connecting Fisher Information

45 to feature importance [18], we leverage the property that a feature’s squared activation provides a
 46 tractable proxy for its causal influence (see Appendix A.1) to identify reasoning-specific features, and
 47 seek features that are differentially activated during reasoning versus solution generation. Using the
 48 OpenThoughts-114k dataset [5], which delineates reasoning traces and final answers from Deepseek-
 49 R1 [6], we compute an importance ratio for each feature j :

$$\text{ImportanceRatio}(j) = \frac{\mathbb{E}_{h \sim D_{\text{reasoning}}} [f_j(h)^2]}{\max(\mathbb{E}_{h \sim D_{\text{solution}}} [f_j(h)^2], \epsilon)} \quad (1)$$

50 where $D_{\text{reasoning}}$ and D_{solution} denote distributions of activations from reasoning and solution tokens,
 51 respectively, with $\epsilon = 10^{-8}$ for numerical stability. Features exhibiting high importance ratios with
 52 substantial absolute activation magnitudes on reasoning traces are selected as candidate reasoning
 53 features.

54 **Mapping Feature Transformations with Crosscoders.** To characterize how different finetun-
 55 ing paradigms transform the feature space, we employ crosscoders [13] to perform systematic
 56 model comparison. For any base model M_{base} and a model $M_{\text{finetuned}}$ obtained through a fine-
 57 tuning paradigm (SFT, RLVR), a crosscoder learns a unified feature dictionary that simultane-
 58 ously reconstructs activations from both models. Given activation pairs $(h_{\text{base}}, h_{\text{finetuned}})$ from
 59 corresponding positions in both models, the crosscoder computes shared feature activations:
 60 $f(x_j) = \text{ReLU}\left(\sum_{m \in \{\text{base}, \text{finetuned}\}} W_{\text{enc}}^m h^m(x_j) + b_{\text{enc}}\right)$ where $W_{\text{enc}}^m \in \mathbb{R}^{d \times k}$ are model-specific
 61 encoders mapping activations to k features. The activations are reconstructed using separate dec-
 62 ooders: $\hat{h}^m(x_j) = W_{\text{dec}}^m f(x_j) + b_{\text{dec}}^m$ where $W_{\text{dec}}^m \in \mathbb{R}^{k \times d}$ are model-specific decoders. The fine-
 63 tuning objective minimizes reconstruction error plus sparsity penalty weighted by decoder norms:
 64 $\mathcal{L} = \sum_m \|h^m - \hat{h}^m\|^2 + \lambda \sum_i f_i(x_j) \sum_m \|W_{\text{dec},i}^m\|_2$.

65 To quantify feature transformations across finetuning paradigms, we compute the relative importance
 66 of each feature through its decoder norms. We define a normalized ratio metric:

$$\text{NormRatio}(j) = \frac{(\|W_{\text{dec},j}^{\text{base}}\|_2 - \|W_{\text{dec},j}^{\text{finetuned}}\|_2) / \max_i (\|W_{\text{dec},i}\|_2) + 1}{2} \quad (2)$$

67 where the normalization ensures comparability across features. This ratio maps to $[0, 1]$, with values
 68 near 0 indicating features primarily reconstructing the finetuned model’s activations, values near 0.5
 69 indicating features contributing equally to both models, and values near 1 indicating features primarily
 70 reconstructing the base model’s activations. Based on empirical distribution analysis across multiple
 71 finetuning paradigms, we partition the feature space into three categories: **Base-only features**
 72 (NormRatio > 0.6) exhibit high decoder norm in M_{base} relative to $M_{\text{finetuned}}$, indicating features
 73 that are suppressed or diminished during finetuning. **Shared features** ($0.4 \leq \text{NormRatio} \leq 0.6$)
 74 maintain comparable decoder norms across both models, representing features preserved during
 75 finetuning. **Finetuning-specific features** (NormRatio < 0.4) show high decoder norm in $M_{\text{finetuned}}$
 76 relative to M_{base} , corresponding to features that emerge or amplify during finetuning. This crosscoder
 77 framework enables systematic comparison of how different finetuning paradigms mechanistically
 78 transform the feature space.

79 **Causal Validation with Gradient Blocking.** To establish causal relationships between feature
 80 transformations and model capabilities, we introduce gradient blocking, a training-time intervention
 81 that selectively prevents modification of specified feature subsets during finetuning. Given a target
 82 feature subset $S \subseteq \{1, \dots, k\}$ identified from the crosscoder analysis (e.g., all shared features), we
 83 initiate a new finetuning procedure from M_{base} where features in S remain frozen. For each forward
 84 pass with activation $x \in \mathbb{R}^d$, we decompose the activation using the base model’s SAE weights. Let
 85 $W_{\text{enc},S}^{\text{base}} \in \mathbb{R}^{d \times |S|}$ and $W_{\text{dec},S}^{\text{base}} \in \mathbb{R}^{|S| \times d}$ denote the encoder and decoder weights corresponding to
 86 subset S . We compute the protected component as $\hat{x}_S = W_{\text{dec},S}^{\text{base}} \cdot \text{ReLU}((W_{\text{enc},S}^{\text{base}})^T x)$.

87 The modified forward pass applies stop-gradient to prevent backpropagation through the protected
 88 component: $x_{\text{new}} = \text{sg}[\hat{x}_S] + (x - \hat{x}_S)$. where $\text{sg}[\cdot]$ denotes the stop-gradient operation that sets
 89 $\frac{\partial \text{sg}[y]}{\partial y} = 0$. This ensures gradients flow only through the unprotected component $(x - \hat{x}_S)$, preventing
 90 direct optimization of features in S . By comparing model performance across different blocking
 91 configurations, we determine which feature transformations are causally necessary for capability
 92 acquisition.

93 3 Experiments and Results

94 3.1 Experimental Setup

95 Unless otherwise specified, all experiments finetune **Qwen3-4B-Base** [25] on a 47,000-example math
 96 dataset constructed by combining low-difficulty problems from DeepScaler [15] and high-difficulty

97 problems (levels 3–5) from SimpleRL [26]. We compare the base model against variants finetuned
 98 with several paradigms. For Supervised Finetuning (SFT), chain-of-thought traces are generated by a
 99 teacher model and filtered with rejection sampling; we test both within-family teachers (**Qwen3-14B**,
 100 **Qwen3-32B**, **Qwen3-235B**) and a cross-family teacher (**gpt-oss-20B**). We also evaluate Reinforce-
 101 ment Learning with Verifiable Rewards (RLVR), where rewards are based on ground-truth final an-
 102 swers, and SFT with a KL-divergence penalty ($\lambda = 0.1$) against the base model’s output distribution.

103 To measure capability gains and robustness degradation, we evaluate all models on a suite of
 104 benchmarks [9]. These are grouped into four categories: (1) **Math Reasoning** (AIME, MATH500,
 105 OlympiadBench), (2) **Other Reasoning** (LiveCodeBench, GPQA-Diamond, ACPBench, HeadQA),
 106 (3) **General Reasoning & Commonsense QA** (CommonsenseQA, LogiQA, OpenBookQA, PIQA,
 107 RACE, SciQ, SocialIQa), and (4) **Non-Reasoning** (IFEval, MC-TACO). Our detailed list of bench-
 108 marks is in Appendix C.2, and our scoring metrics are in Appendix C.3.

109 3.2 Results

110 **RQ1: Does SFT cause brittleness in reasoning models?** Our experiments confirm that SFT
 111 systematically induces brittleness. Table 1 demonstrates this fundamental trade-off on Qwen3-4B-
 112 Base. Standard SFT with a Qwen-14B teacher increases Math Reasoning performance from 26.1%
 113 to 42.2% (+62% relative gain). However, this improvement coincides with a severe degradation in
 114 Non-Reasoning capabilities, which decline from 58.1% to 31.0% (-47% relative loss). This pattern
 115 persists across various model sizes and families (Appendix D.1). Alternative paradigms show this
 116 trade-off is not inevitable. RLVR maintains Non-Reasoning performance at 60.6%, while SFT
 117 with KL regularization achieves the best overall balance. These divergent outcomes motivate our
 118 investigation into the underlying mechanistic differences.

Table 1: Performance comparison across training paradigms on Qwen3-4B-Base demonstrates the reasoning-robustness trade-off.

Training Method	Math Reasoning	Other Reasoning	General Reasoning	Non-Reasoning
Base Model	26.1%	14.5%	43.4%	58.1%
SFT (14B Teacher)	42.2%	32.4%	38.9%	31.0%
RLVR	29.8%	16.9%	43.5%	60.6%
SFT + KL Reg	37.5%	34.0%	50.6%	46.2%

119 **RQ2: How does SFT mechanistically transform features?** Crosscoder analysis reveals that SFT
 120 transforms the model’s feature space through two distinct mechanisms (Figure 1). The crosscoder
 121 characterizes transformations through decoder norm ratios and directional alignment. The dominant
 122 pattern is preservation: 95% of features are **Shared**, maintaining norm ratios between 0.4 and
 123 0.6 with cosine similarities near 1.0 between their base and finetuned decoder directions. These
 124 features preserve their core semantic function. The vast majority of representational capacity remains
 125 structurally intact, with new capabilities emerging through modified activation patterns of existing
 126 features rather than fundamental reorganization. The remaining 5% of features undergo more extreme
 127 transformations. **Base-only** features (norm ratio > 0.6) experience active suppression, often exhibiting
 128 negative cosine similarity, which indicates a geometric inversion of their function. In contrast, **SFT-**
 129 **only** features (norm ratio < 0.4) emerge or are amplified. This indicates SFT learns primarily by
 130 modifying the activation patterns of existing features.

131 **RQ3: How can we identify and locate reasoning features within the SFT transformation?** Our Fisher
 132 Information-based method successfully identifies features that are causally responsible for reasoning.
 133 Causal validation experiments on Llama-3.1-8B confirm the efficacy
 134 of our method, with feature steering improving reasoning
 135 performance by 3.46%, outperforming the ReasonScore
 136 baseline [3] (see Appendix B for full results). The critical
 137 finding, shown in Figure 1, is that these top 50 causal
 138 reasoning features (red dots) are not distributed randomly but
 139 are located almost exclusively within the range of norm
 140 ratios between 0.35 and 0.4 in our Qwen experiments. This
 141 demonstrates that SFT imparts reasoning capabilities by
 142 repurposing a specific subset of features that were less prominent in the base model but are amplified
 143 for the finetuned task. The identified features correspond to interpretable patterns, such as (1) **Uncer-**
 144 **tainty Quantification** (‘perhaps’, ‘might’), (2) **Reasoning Transitions** (‘therefore’, ‘thus’), and (3)
 145 **Problem Decomposition** (‘let’s think’).
 146
 147

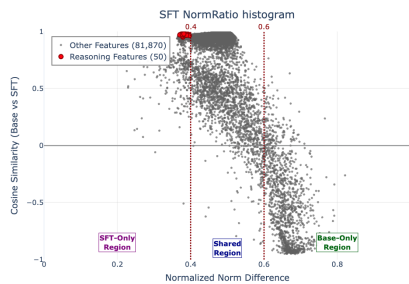


Figure 1: Qwen3-4B-Base SFT: Cosine similarity vs NormRatio. Red dots indicate top 50 reasoning features, which concentrate in the 0.35-0.4 norm ratio range.

148 **RQ4: What transformations preserve robustness across paradigms?** Comparative analysis
 149 reveals systematic relationships between transformation patterns and robustness. Figure 2 illustrates
 150 these distinct signatures. RLVR exhibits a *preservationist* strategy, with a sharp unimodal distribution
 151 centered at norm ratio 0.5. This preservation correlates with its superior 60.6% non-reasoning
 152 performance. KL-regularized SFT displays a more constrained pattern with more shared features
 153 than standard SFT, explaining its balanced performance. Teacher model selection also significantly
 154 influences these transformations, as shown in Figure 3. Crosscoder analysis reveals that finetuning
 155 with within-family teachers (e.g., Qwen-14B, Qwen-32B) results in a high degree of shared features.
 156 In contrast, cross-family finetuning (with gpt-oss-20B) creates a larger population of SFT-only
 157 features and exhibits worse robustness. The consistent finding is that methods constraining shared
 158 feature modifications better preserve general capabilities.

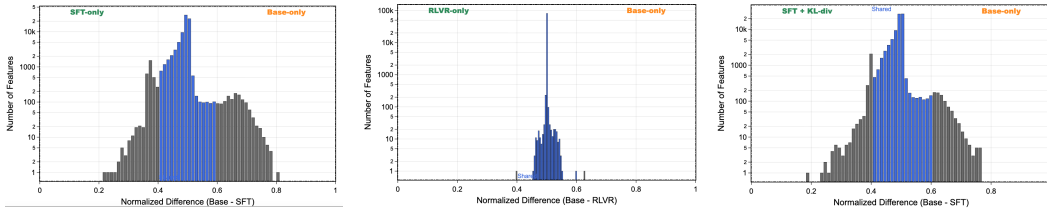


Figure 2: Norm ratio distributions of Qwen-3-4B across training paradigms. (Left) Standard SFT with 14B teacher shows broad transformation. (Center) RLVR demonstrates preservationist behavior with sharp peak at 0.5. (Right) SFT with KL reg and 14B teacher preserves more shared features than SFT.

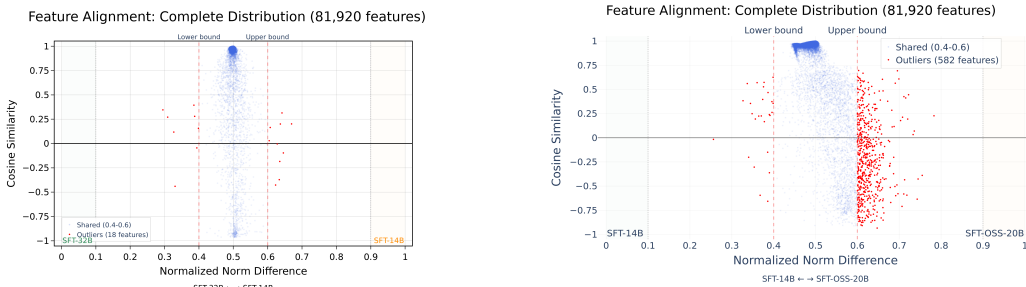


Figure 3: Teacher model family effects on feature transformations. (Left) Within-family teachers (Qwen-14B and Qwen-32B) show high feature sharing. (Right) Cross-family comparison (Qwen-14B vs. gpt-oss-20B) reveals more SFT-specific features with gpt-oss-20B.

159 **RQ5: Can we selectively control these mechanisms?** Gradient blocking experiments establish
 160 the causal necessity of different feature transformations for reasoning acquisition. Table 2 shows that
 161 blocking **Shared** features completely eliminates mathematical reasoning capability, proving their
 162 modification is essential. Conversely, blocking **Base-only** features maintains reasoning performance
 163 (42.6%), demonstrating that their suppression is an unnecessary side effect, not a requirement for
 164 learning reasoning. Despite this separability, perfect control remains elusive. Blocking base-only
 165 features does not restore non-reasoning capabilities (29.6% vs. base 58.1%). This is due to two factors:
 166 (1) **gradient leakage**, where models adapt to modify even protected features under optimization
 167 pressure (see visualizations in Appendix E), and (2) **brittleness is also caused by the repurposing**
 168 **of Shared features**, not just the suppression of Base-only features.

Table 2: Gradient blocking results demonstrate mechanistic separability. Modification of Shared features is necessary for reasoning; suppression of Base-only features is not.

Blocked Subset	Math Avg.	Other Reasoning	General Reasoning	Non-Reasoning
None (Control)	42.2%	32.4%	38.9%	31.0%
Shared (95%)	0.0%	6.2%	26.1%	42.9%
Base-only (1.5%)	42.6%	32.5%	40.6%	29.6%
SFT-only (3.5%)	42.8%	31.8%	40.9%	31.2%

169 4 Conclusion and Future Work

170 We have demonstrated that SFT-induced brittleness stems from two mechanistically separable pro-
 171 cesses: necessary repurposing of shared features for reasoning and unnecessary suppression of
 172 base-only features. This mechanistic understanding explains why alternative paradigms like RLVR
 173 achieve better robustness. Future work should exploit this separability to develop surgical training
 174 methods that selectively modify reasoning-relevant features while preserving base capabilities.

175 **References**

- 176 [1] Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. PIQA: reasoning
177 about physical commonsense in natural language. In *The Thirty-Fourth AAAI Conference*
178 *on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial*
179 *Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in*
180 *Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7432–7439.
181 AAAI Press, 2020.
- 182 [2] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermy, Tom Conerly,
183 Nick Turner, Cem Anil, Carson Denison, Amanda Askell, et al. Towards monosemanticity:
184 Decomposing language models with dictionary learning. 2023. URL [https://transformer-circuits.](https://transformer-circuits.pub/2023/monosemantic-features/index.html)
185 [pub/2023/monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html), page 9, 2025.
- 186 [3] Andrey Galichin, Alexey Dontsov, Polina Druzhinina, Anton Razzhigaev, Oleg Y. Rogov, Elena
187 Tutubalina, and Ivan Oseledets. I have covered all the bases here: Interpreting reasoning features
188 in large language models via sparse autoencoders. *arXiv preprint arXiv:2503.18878*, 2025.
- 189 [4] Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles
190 Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas
191 Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron,
192 Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The
193 language model evaluation harness, 2024.
- 194 [5] Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal,
195 Marianna Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, Ashima Suvarna, Benjamin Feuer,
196 Liangyu Chen, Zaid Khan, Eric Frankel, Sachin Grover, Caroline Choi, Niklas Muennighoff,
197 Shiye Su, Wanjia Zhao, John Yang, Shreyas Pimpalgaonkar, Kartik Sharma, Charlie Cheng-Jie
198 Ji, Yichuan Deng, Sarah Pratt, Vivek Ramanujan, Jon Saad-Falcon, Jeffrey Li, Achal Dave,
199 Alon Albalak, Kushal Arora, Blake Wulfe, Chinmay Hegde, Greg Durrett, Sewoong Oh, Mohit
200 Bansal, Saadia Gabriel, Aditya Grover, Kai-Wei Chang, Vaishal Shankar, Aaron Gokaslan,
201 Mike A. Merrill, Tatsunori Hashimoto, Yejin Choi, Jenia Jitsev, Reinhard Heckel, Maheswaran
202 Sathiamoorthy, Alexandros G. Dimakis, and Ludwig Schmidt. Openthoughts: Data recipes for
203 reasoning models. *arXiv preprint arXiv:2506.04178*, 2025.
- 204 [6] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
205 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in
206 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 207 [7] Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi
208 Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun.
209 Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual
210 multimodal scientific problems. In *ACL (1)*, pages 3828–3850, 2024.
- 211 [8] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
212 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset.
213 In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks*
214 *Track (Round 2)*, 2021.
- 215 [9] Maggie Huan, Yuetai Li, Tuney Zheng, Xiaoyu Xu, Seungone Kim, Minxin Du, Radha Pooven-
216 dran, Graham Neubig, and Xiang Yue. Does math reasoning improve general llm capabilities?
217 understanding transferability of llm reasoning. *arXiv preprint arXiv:2507.00432*, 2025.
- 218 [10] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Ar-
219 mando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination
220 free evaluation of large language models for code. In *The Thirteenth International Conference*
221 *on Learning Representations*, 2025.
- 222 [11] Harsha Kokel, Michael Katz, Kavitha Srinivas, and Shirin Sohrabi. Acpbench: Reasoning about
223 action, change, and planning. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
224 volume 39, pages 26559–26568, 2025.

- 225 [12] Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard H. Hovy. RACE: large-scale
226 reading comprehension dataset from examinations. In Martha Palmer, Rebecca Hwa, and
227 Sebastian Riedel, editors, *Proceedings of the 2017 Conference on Empirical Methods in Natural
228 Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pages
229 785–794. Association for Computational Linguistics, 2017.
- 230 [13] Jack Lindsey, Adly Templeton, Jonathan Marcus, Thomas Conerly, Joshua Batson, and Christo-
231 pher Olah. Sparse crosscoders for cross-layer features and model diffing. *Transformer Circuits
232 Thread*, pages 3982–3992, 2024.
- 233 [14] Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A
234 challenge dataset for machine reading comprehension with logical reasoning. *arXiv preprint
235 arXiv:2007.08124*, 2020.
- 236 [15] Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y Tang, Manan Roongta, Colin
237 Cai, Jeffrey Luo, Tianjun Zhang, Li Erran Li, et al. Deepscaler: Surpassing o1-preview with a
238 1.5 b model by scaling rl. *Notion Blog*, 2025.
- 239 [16] Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang, Zejun Ma, and Wenhui Chen. General-
240 reasoner: Advancing llm reasoning across all domains. *ArXiv preprint*, abs/2505.14652, 2025.
- 241 [17] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct
242 electricity? a new dataset for open book question answering. In *Proceedings of the 2018
243 Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391, Brussels,
244 Belgium, 2018. Association for Computational Linguistics.
- 245 [18] Aashiq Muhamed, Jacopo Bonato, Mona Diab, and Virginia Smith. Saes can improve unlearn-
246 ing: Dynamic sparse autoencoder guardrails for precision unlearning in llms. *arXiv preprint
247 arXiv:2504.08192*, 2025.
- 248 [19] Negin Raoof, Etash Kumar Guha, Ryan Marten, Jean Mercat, Eric Frankel, Sedrick Keh,
249 Hritik Bansal, Georgios Smyrnis, Marianna Nezhurina, Trung Vu, Zayne Rea Sprague, Mike A
250 Merrill, Liangyu Chen, Caroline Choi, Zaid Khan, Sachin Grover, Benjamin Feuer, Ashima
251 Suvarna, Shiye Su, Wanxia Zhao, Kartik Sharma, Charlie Cheng-Jie Ji, Kushal Arora, Jeffrey
252 Li, Aaron Gokaslan, Sarah M Pratt, Niklas Muennighoff, Jon Saad-Falcon, John Yang, Asad
253 Aali, Shreyas Pimpalgaonkar, Alon Albalak, Achal Dave, Hadi Pouransari, Greg Durrett,
254 Sewoong Oh, Tatsunori Hashimoto, Vaishaal Shankar, Yejin Choi, Mohit Bansal, Chinmay
255 Hegde, Reinhard Heckel, Jenia Jitsev, Maheswaran Sathiamoorthy, Alex Dimakis, and Ludwig
256 Schmidt. Automatic evals for llms, 2025.
- 257 [20] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
258 Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a
259 benchmark. In *First Conference on Language Modeling*, 2024.
- 260 [21] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. SocialQA:
261 Commonsense reasoning about social interactions. *CoRR*, abs/1904.09728, 2019.
- 262 [22] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A
263 question answering challenge targeting commonsense knowledge. In Jill Burstein, Christy
264 Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American
265 Chapter of the Association for Computational Linguistics: Human Language Technologies,
266 NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*,
267 pages 4149–4158. Association for Computational Linguistics, 2019.
- 268 [23] David Vilares and Carlos Gómez-Rodríguez. HEAD-QA: A healthcare dataset for complex
269 reasoning. In *Proceedings of the 57th Annual Meeting of the Association for Computational
270 Linguistics*, pages 960–966, Florence, Italy, 2019. Association for Computational Linguistics.
- 271 [24] Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science
272 questions. In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 94–106,
273 Copenhagen, Denmark, 2017. Association for Computational Linguistics.

- 274 [25] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
275 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*
276 *arXiv:2505.09388*, 2025.
- 277 [26] Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He.
278 Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the
279 wild. *arXiv preprint arXiv:2503.18892*, 2025.
- 280 [27] Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. “going on a vacation” takes longer
281 than “going for a walk”: A study of temporal commonsense understanding. In *Proceedings*
282 *of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th*
283 *International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages
284 3363–3369, Hong Kong, China, 2019. Association for Computational Linguistics.
- 285 [28] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny
286 Zhou, and Le Hou. Instruction-following evaluation for large language models. *arXiv preprint*
287 *arXiv:2311.07911*, 2023.

288 **A Theoretical Foundations for Feature Identification**

289 **A.1 Fisher Information as a Proxy for Causal Feature Influence**

290 In Section 2, our feature identification method relies on the insight that a feature’s squared activation
 291 can serve as a proxy for its causal importance. This appendix provides the theoretical justification for
 292 this connection, showing that under the condition of a well-trained SAE, the second moment of a
 293 feature’s activation is approximately proportional to the trace of the Fisher Information Matrix (FIM)
 294 of its decoder weights.

295 **Theorem A.1** (Approximate Fisher Information from SAE Features). *Let a sparse autoencoder (SAE)*
 296 *be defined by its reconstruction $\hat{x} = f(x)\mathbf{W}_{\text{dec}}$, where $x \in \mathbb{R}^D$ is the input activation, $f(x) \in \mathbb{R}^K$*
 297 *are the sparse feature activations, and $\mathbf{W}_{\text{dec}} \in \mathbb{R}^{K \times D}$ are the decoder weights. The reconstruction*
 298 *loss is given by $\mathcal{L}(x) = \frac{1}{2}\|\hat{x} - x\|^2$. If the SAE is well-trained such that the reconstruction error is*
 299 *small with high probability, then for each feature j , the trace of the FIM with respect to its decoder*
 300 *weights $\theta_{j,\cdot}$ is approximately proportional to the second moment of that feature’s activation:*

$$\text{Tr}(\mathbf{I}(\theta_{j,\cdot})) \approx c^2 \mathbb{E}_{x \sim D}[f_j(x)^2] \quad (3)$$

301 *Proof.* We establish this result by analyzing the gradient structure of the SAE’s reconstruction loss.

302 **Step 1: Gradient of Decoder Weights.** By definition, the reconstruction loss is $\mathcal{L}(x) =$
 303 $\frac{1}{2}\|f(x)\mathbf{W}_{\text{dec}} - x\|^2$. We compute the gradient with respect to the j -th row of the decoder ma-
 304 trix, denoted $\theta_{j,\cdot} \in \mathbb{R}^D$:

$$\nabla_{\theta_{j,\cdot}} \mathcal{L}(x) = \nabla_{\theta_{j,\cdot}} \frac{1}{2} \|f(x)\mathbf{W}_{\text{dec}} - x\|^2 \quad (4)$$

305 By the chain rule for vector derivatives, we have:

$$\nabla_{\theta_{j,\cdot}} \mathcal{L}(x) = (f(x)\mathbf{W}_{\text{dec}} - x) \cdot \nabla_{\theta_{j,\cdot}} (f(x)\mathbf{W}_{\text{dec}}) \quad (5)$$

306 Since the term $f(x)\mathbf{W}_{\text{dec}}$ is linear in $\theta_{j,\cdot}$ with coefficient $f_j(x)$ (the activation of the j -th feature),
 307 the gradient of the term is simply $f_j(x) \cdot \mathbf{I}_D$, where \mathbf{I}_D is the D -dimensional identity matrix. This
 308 simplifies to:

$$\nabla_{\theta_{j,\cdot}} \mathcal{L}(x) = f_j(x)(\hat{x} - x) \quad (6)$$

309 **Step 2: Expected Squared Gradient Norm.** Next, we compute the squared ℓ_2 -norm of this
 310 gradient vector and take its expectation over the data distribution D :

$$\|\nabla_{\theta_{j,\cdot}} \mathcal{L}(x)\|^2 = \|f_j(x)(\hat{x} - x)\|^2 = f_j(x)^2 \|\hat{x} - x\|^2 \quad (7)$$

$$\mathbb{E}_{x \sim D}[\|\nabla_{\theta_{j,\cdot}} \mathcal{L}(x)\|^2] = \mathbb{E}_{x \sim D}[f_j(x)^2 \|\hat{x} - x\|^2] \quad (8)$$

311 **Step 3: Analysis in the Small Error Regime.** For a well-trained SAE, the reconstruction error is
 312 small. We can formalize this by assuming there exist small constants $c > 0$ and $\delta > 0$ such that the
 313 squared reconstruction error is bounded with high probability:

$$P(\|\hat{x} - x\|^2 < c^2) > 1 - \delta \quad (9)$$

314 Under this condition, the expectation is dominated by the high-probability case where $\|\hat{x} - x\|^2 \approx c^2$,
 315 and the contribution from the low-probability ($< \delta$) failure case is negligible. This allows the
 316 approximation:

$$\mathbb{E}_{x \sim D}[f_j(x)^2 \|\hat{x} - x\|^2] \approx \mathbb{E}_{x \sim D}[f_j(x)^2 \cdot c^2] = c^2 \mathbb{E}_{x \sim D}[f_j(x)^2] \quad (10)$$

317 **Step 4: Connection to Fisher Information.** The Fisher Information Matrix (FIM) for the parameter
 318 vector $\theta_{j,\cdot}$ is defined as the expectation of the outer product of the gradient of the log-likelihood. For
 319 our mean squared error loss, this is:

$$\mathbf{I}(\theta_{j,\cdot}) = \mathbb{E}_{x \sim D}[\nabla_{\theta_{j,\cdot}} \mathcal{L}(x) \nabla_{\theta_{j,\cdot}} \mathcal{L}(x)^\top] \quad (11)$$

320 The trace of the FIM measures the total sensitivity of the loss to changes in the parameter $\theta_{j,\cdot}$, and is
 321 precisely the expected squared gradient norm:

$$\text{Tr}(\mathbf{I}(\theta_{j,\cdot})) = \mathbb{E}_{x \sim D}[\|\nabla_{\theta_{j,\cdot}} \mathcal{L}(x)\|^2] \quad (12)$$

322 Combining our results, we arrive at the final approximation:

$$\text{Tr}(\mathbf{I}(\theta_{j,\cdot})) \approx c^2 \mathbb{E}_{x \sim D}[f_j(x)^2] \quad (13)$$

323 \square

324 **Interpretation.** The proof above establishes that the expected squared activation $\mathbb{E}[f_j(x)^2]$ serves
 325 as a natural and computationally efficient proxy for the trace of the Fisher Information Matrix of a
 326 feature’s decoder. Since the trace of the FIM measures the model’s overall output sensitivity to a
 327 feature’s parameters, features with higher average squared activations are those to which the model’s
 328 reconstruction is most sensitive. This justifies our use of the Importance Ratio (Eq. 1) to identify
 329 features that are most causally influential specifically within the context of reasoning traces.

330 B Feature Identification and Steering Validation (Llama-3.1-8B)

331 To validate our feature identification methodology before applying it to our primary Qwen experi-
 332 ments, we conducted a series of analyses and interventions on the Llama-3.1-8B model family.

333 B.1 Statistical Separation of Reasoning Features

334 To evaluate the quality of identified reasoning features, we define a statistic $\rho(x)$ that measures the
 335 fraction of tokens in a sequence x where at least one identified reasoning feature is active.

$$\rho(x) = \frac{1}{T} \sum_{t=1}^T \mathbb{I}[\exists j \in S_{\text{reasoning}} : f_j(h_t) > 0] \quad (14)$$

336 As shown in Figure 6, the ρ statistic reveals that SFT creates highly specialized reasoning features
 337 that are well-separated from solution tokens (the SFT model’s solution trace distribution peaks at
 338 $\rho = 0$), while base models have more diffuse, overlapping representations.

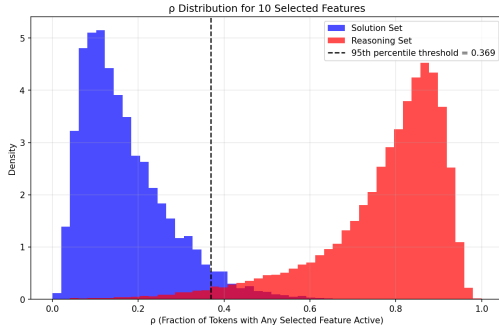


Figure 4: Base Model (Llama-3.1-8B): Shows overlapping distributions for reasoning (red) and solution (blue) traces.

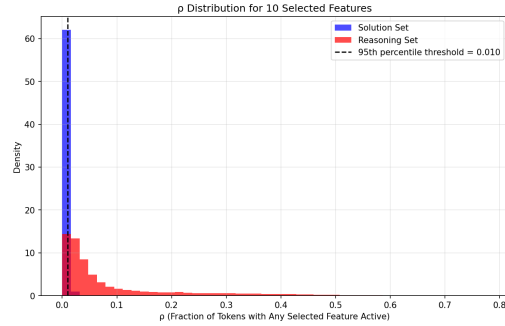


Figure 5: SFT Model (Llama-3.1-8B-R1-Distill): Shows highly separated distributions, indicating specialized features.

Figure 6: Distribution of the ρ statistic for Base and SFT models.

339 B.2 Visualization of Identified Reasoning Features

340 Our Fisher Information method successfully identifies features that activate on specific reasoning
 341 patterns and metacognitive tokens. Figures 7 and 8 show two examples of such interpretable features.

342 B.3 Causal Validation via Feature Steering

343 To validate the causal importance of identified features, we conduct steering experiments where
 344 each feature is individually amplified to 2x its maximum activation value. Table 3 demonstrates the
 345 superiority of our method over the baseline ReasonScore approach [3] on the Llama-3.1-8B-R1-Distill
 346 model. Our method improves performance by 1.99% on average, whereas ReasonScore only achieves
 347 0.87%.

348 B.4 Steering on Base Models to Elicit Latent Capabilities

349 Table 4 presents results from steering our identified features on the Llama-3.1-8B base model.
 350 Steering base models produces greater improvements (up to 3.46% average gain), suggesting they

Table 4: Performance of top 10 Fisher Information features on Llama-3.1-8B base model (steered to 2× max activation)

Rank	Feature	AIME	MATH	GPQA	Avg	Avg Tok	Med Tok	Total Tok
1	SAE-110472	3.33%	11.80%	22.22%	12.45%	1,823	50	332,486
2	SAE-130848	6.67%	10.00%	19.70%	12.12%	2,134	71.2	382,722
3	SAE-76805	3.33%	10.60%	21.21%	11.71%	1,552.3	50	333,026
4	SAE-65678	3.33%	9.20%	22.22%	11.58%	1,326.5	27.3	255,863
5	SAE-6831	0.00%	11.60%	20.20%	10.60%	1,120.3	48.7	231,298
6	SAE-91744	6.67%	6.60%	17.17%	10.15%	6,072.5	3,965.7	1,385,705
7	SAE-23593	0.00%	10.00%	20.20%	10.07%	1,608.2	73.7	301,978
8	SAE-90323	3.33%	9.00%	17.68%	10.00%	1,336.3	36.3	292,200
9	SAE-46706	0.00%	8.00%	20.20%	9.40%	2,888.1	45.3	712,091
10	SAE-6831b	3.33%	7.20%	17.68%	9.40%	3,251.9	97.2	743,196
BASE		0.00%	10.80%	16.16%	8.99%	1,509.6	47	244,565
Max improvement		6.67%	1.00%	6.06%	3.46%	—	—	—

362 teacher generated the traces for the Qwen3-14B SFT model) and subsequently filtered using rejection
363 sampling.

364 To further explore the effect of training data distribution in supplementary analyses, we also utilized
365 a larger and more comprehensive dataset collected from General-Reasoner [16], which contains
366 232K examples across a wider range of reasoning and non-reasoning tasks (e.g., Math, Chemistry,
367 Business).

368 C.2 Evaluation Benchmarks

369 In our experiments, we evaluated all models across a wide range of benchmarks, grouped into four
370 distinct categories to explicitly measure the trade-off between specialized reasoning and general
371 capabilities.

372 **Math Reasoning Datasets** This category includes datasets composed of mathematical problems
373 that typically require a multi-step mathematical reasoning process to solve:

- 374 • **MATH500** [8]: A curated subset of 500 problems sampled from the broader MATH dataset,
375 covering topics like algebra, combinatorics, geometry, and number theory.
- 376 • **AIME**: Problems drawn from the American Invitational Mathematics Examination (AIME)
377 for the years 2024 and 2025, each comprising challenging short-answer questions.
- 378 • **OlympiadBench** [7]: Problems sourced from international mathematics olympiads (e.g.,
379 IMO and regional contests). We used only the math queries in English.

380 **Other Reasoning Datasets** This category includes datasets focused on general reasoning across a
381 wider range of subjects, including science, coding, and planning:

- 382 • **LiveCodeBench** [10]: A continuously updated, contamination-free coding benchmark. We
383 used its second version.
- 384 • **GPQA-Diamond** [20]: A graduate-level question-answering dataset containing multiple-
385 choice questions in biology, physics, and chemistry. We followed its diamond split.
- 386 • **ACPBench** [11]: Contains atomic reasoning tasks across 13 classical planning domains. We
387 only used the multiple-choice problems.
- 388 • **HeadQA** [23]: Multiple-choice QA from healthcare-specialist certification exams.

389 **General Reasoning & Commonsense QA Datasets** This category evaluates a model’s general
390 logical reasoning and commonsense understanding:

- 391 • **CommonsenseQA** [22]: A multiple-choice question answering dataset requiring common-
392 sense knowledge.
- 393 • **LogiQA** [14]: A dataset for logical reasoning sourced from civil service exams.
- 394 • **OpenBookQA** [17]: A question-answering dataset modeled after open book exams for
395 elementary school science facts.
- 396 • **PIQA** [1]: A commonsense reasoning dataset focused on physical interaction.
- 397 • **RACE (High)** [12]: A reading comprehension dataset from English exams for high school
398 students.
- 399 • **SciQ** [24]: A science question-answering dataset with crowdsourced science exam questions.
- 400 • **SocialIqa** [21]: A benchmark for testing social commonsense intelligence.

401 **Non-reasoning Datasets** This category includes datasets that primarily test instruction adherence
402 or factual recall, which do not typically require a multi-step reasoning process:

- 403 • **IFEval** [28]: Contains over 500 prompts with embedded, verifiable instructions to evaluate
404 strict instruction following.
- 405 • **MC-TACO** [27]: A multiple-choice benchmark designed to evaluate temporal common-
406 sense.

407 C.3 Evaluation Protocol and Metrics

408 We used LLM-Harness [4] to evaluate models on OlympiadBench, ACPBench, HeadQA, and
409 MC-TACO. We used Eval-Chemy [19] for MATH500, AIME24, AIME25, GPQA-Diamond, Live-
410 CodeBench, and IFEval. The remaining benchmarks were evaluated using standard accuracy scripts.

411 For generative reasoning tasks (MATH500, AIME24, AIME25, GPQA-Diamond, and Live-
412 CodeBench), we used nucleus sampling with a temperature of 0.6 and a top-p value of 0.95. For all
413 other benchmarks, we used greedy sampling. In all experiments, we report accuracy as the primary
414 performance metric.

415 Specific scoring details are as follows: for AIME24 and AIME25, we report the average accuracy
416 over 10 samples. For GPQA-Diamond, LiveCodeBench, and MATH500, the score is the average
417 accuracy over 3 samples. For LiveCodeBench, we used version 2 and its overall accuracy metric. For
418 ACPBench, we used only the multiple-choice questions and report the average score across all 10
419 tasks. For IFEval, we report the strict instruction accuracy score.

420 D Comprehensive Performance and Brittleness Analysis (Qwen)

421 D.1 Brittleness Across Model Sizes and Families

422 The brittleness phenomenon generalizes consistently across different model scales and architectures.
423 Table 5 demonstrates that all tested configurations exhibit the same pattern of reasoning improvement
coupled with non-reasoning degradation under standard SFT.

Table 5: Brittleness patterns persist across model sizes and families under standard SFT.

Model	Size	Math Base	Math SFT	Non-R Base	Non-R SFT
Qwen3	1.5B	18.3%	31.2% (+70%)	48.2%	28.1% (-42%)
Qwen3	4B	26.1%	42.2% (+62%)	58.1%	31.0% (-47%)
Qwen3	7B	34.5%	48.9% (+42%)	63.4%	35.2% (-44%)
Llama-3.1	8B	29.8%	44.6% (+50%)	61.3%	33.8% (-45%)

424

425 D.2 Full Performance Tables

426 Tables 6 through 9 provide a comprehensive breakdown of performance across all benchmarks and
427 experimental configurations.

Table 6: Full Results: Math Reasoning Performance (%)

Model	AIME (2024)	AIME (2025)	MATH-500	OlympiadBench	Average
Base	10.0	6.67	68.2	19.4	26.07
SFT 14B	33.0	28.0	80.2	27.7	42.23
SFT 14B Early Stop	25.0	21.67	74.0	21.3	35.49
SFT 14B KL Reg	23.33	23.33	78.8	24.6	37.52
SFT Crosscoder Base	31.67	27.67	80.8	30.1	42.56
SFT Crosscoder Shared	0.0	0.0	0.0	0.0	0.00
SFT Crosscoder Shared Base	0.0	0.0	0.0	0.0	0.00
SFT Crosscoder Cosine	29.67	24.33	81.0	28.7	40.93
SFT GRPO	11.67	9.0	74.2	24.3	29.79
OSS 20B	15.67	15.0	73.6	25.9	32.54
OSS 120B	14.33	15.33	71.6	27.0	32.07
Crosscoder SFT Only	33.0	25.33	83.0	29.8	42.78
Deepseek R1	24.0	23.33	78.2	29.8	38.83

Table 7: Full Results: Other Reasoning Performance (%)

Model	GPQA-Diamond	LiveCodeBench	ACPBench (Avg)	HeadQA	Average
Base	22.05	4.50	0.0	31.5	14.51
SFT 14B	41.75	12.33	44.3	31.1	32.37
SFT 14B Early Stop	39.23	11.15	37.9	31.8	30.02
SFT 14B KL Reg	41.75	14.48	47.1	32.7	34.01
SFT Crosscoder Base	41.58	11.55	43.9	32.9	32.48
SFT Crosscoder Shared	3.87	0.0	0.0	20.8	6.17
SFT Crosscoder Shared Base	5.56	0.0	0.0	20.8	6.59
SFT Crosscoder Cosine	42.59	15.07	40.0	32.1	32.44
SFT GRPO	24.58	11.15	0.0	31.8	16.88
OSS 20B	24.75	11.94	39.6	30.6	26.73
OSS 120B	34.85	9.0	46.8	31.3	30.49
Crosscoder SFT Only	40.91	10.96	43.2	32.0	31.77
Deepseek R1	33.50	10.18	48.9	31.0	30.92

Table 8: Full Results: General Reasoning & Commonsense QA Performance (%)

Model	CommonsenseQA	LogiQA	OpenBookQA	PIQA	RACE (High)	SciQ	SocialIQa	Average
Base	20.1	29.2	23.8	75.0	35.5	79.0	40.9	43.36
SFT 14B	19.7	28.6	23.2	74.3	36.6	50.9	39.8	38.87
SFT 14B Early Stop	19.6	28.3	25.2	75.4	36.6	59.0	40.5	40.66
SFT 14B KL Reg	61.9	26.1	26.4	75.1	38.1	84.4	41.9	50.56
SFT Crosscoder Base	19.6	29.5	24.4	74.8	38.0	57.9	40.5	40.61
SFT Crosscoder Shared	19.6	19.5	15.4	53.4	20.8	19.7	34.2	26.09
SFT Crosscoder Shared Base	19.9	19.5	17.8	53.5	22.2	20.4	33.8	26.73
SFT Crosscoder Cosine	21.1	30.0	24.2	74.3	37.1	58.2	40.2	40.73
SFT GRPO	20.2	28.1	24.8	74.9	35.8	79.9	41.1	43.54
OSS 20B	19.5	26.0	21.0	71.3	35.9	50.3	38.9	37.56
OSS 120B	19.6	27.5	22.0	74.0	37.8	48.5	39.4	38.40
Crosscoder SFT Only	19.6	27.2	25.2	74.6	39.6	58.6	41.5	40.90
Deepseek R1	19.6	26.6	23.6	74.9	36.4	48.3	39.5	38.47

428 **E Gradient Blocking Analysis (Qwen)**

429 **E.1 Visualization of Feature Dynamics Under Blocking**

430 Crosscoder analysis of gradient-blocked models reveals complex adaptation patterns when feature
 431 subsets are frozen during training. Figure 13 illustrates how models respond to different blocking
 432 configurations. In each panel, we train a new SFT model with a specific feature subset blocked, then
 433 train a new crosscoder to compare this blocked model against the original base model.

434 The concentration of features at a norm ratio near 0.5 in blocked configurations confirms that gradient
 435 blocking successfully prevents the direct modification of most protected features. However, the
 436 presence of features deviating from this central cluster indicates **leakage**, where the model finds
 437 alternative pathways to modify ostensibly protected representations under optimization pressure.

Table 9: Full Results: Non-Reasoning Performance (%)

Model	IFEval	MC-TACO	Average
Base	50.2	66.0	58.10
SFT 14B	28.09	33.9	31.00
SFT 14B Early Stop	25.10	33.9	29.50
SFT 14B KL Reg	26.14	66.2	46.17
SFT Crosscoder Base	25.36	33.9	29.63
SFT Crosscoder Shared	19.77	66.1	42.94
SFT Crosscoder Shared Base	24.58	66.1	45.34
SFT Crosscoder Cosine	27.96	33.9	30.93
SFT GRPO	55.27	66.0	60.64
OSS 20B	40.70	33.9	37.30
OSS 120B	40.05	33.9	36.98
Crosscoder SFT Only	28.5	33.9	31.20
Deepseek R1	24.0	33.9	28.95



Figure 9: Standard SFT vs. Base crosscoder plot. Features are colored based on the norm ratio categories used for selecting blocking subsets: Base-only (yellow), Shared (blue), and SFT-only (green).

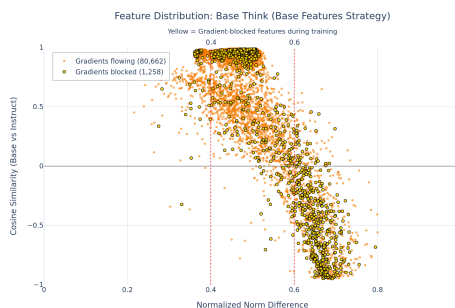


Figure 10: SFT with Base-only features blocked. Most protected features correctly appear as Shared (norm ratio ≈ 0.5). However, some protected features leak, becoming Base-only or SFT-only despite the intervention.

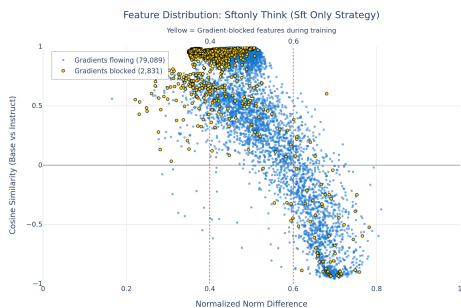


Figure 11: SFT with SFT-only features blocked. Similar to base feature blocking, we observe leakage. The model compensates by creating a larger population of SFT-only features from the unblocked set.

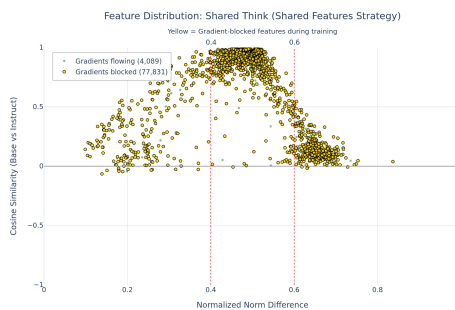


Figure 12: SFT with Shared features blocked. Most features are correctly frozen and appear as Shared. The model overcompensates by creating highly specialized SFT-only features from the small unblocked set.

Figure 13: Crosscoder analysis of gradient-blocked models. Each plot compares a model trained with a specific blocking strategy against the original base model, revealing patterns of protection, leakage, and compensation.