

Neuron-Aware Data Selection for Annotation-Free LLM Self-Distillation

Anonymous Author(s)
 Affiliation
 Address
 email

Abstract

1 Post-training large language models (LLMs) without real-world interaction feedback or human-labeled supervision remains challenging, particularly in specialized
 2 domains where expert annotations are costly to obtain. Recent annotation-free self-
 3 evolution methods address this by using the model’s own outputs as supervision
 4 signals, constructing a teacher via additional context and aggregating predictions
 5 across multiple rollouts through majority voting to produce pseudo-labels. How-
 6 ever, these approaches are not without drawbacks: SFT- and GRPO-based variants
 7 suffer out-of-domain performance degradation, while reward-based on-policy RL
 8 inflates calibration error. In this paper, we propose NEURON ON-POLICY SELF-
 9 DISTILLATION (NEURON-OPSD), a data-centric framework for annotation-free
 10 self-distillation that leverages internal neuron activations to guide both training-
 11 data selection and teacher context construction. The model is then trained via
 12 on-policy distillation from the teacher distribution, requiring no ground-truth labels
 13 at any stage. Across specialized-domain benchmarks, NEURON-OPSD improves
 14 in-domain task performance while preserving cross-domain generalization and mit-
 15 igating calibration collapse over prior annotation-free baselines. This framework is
 16 particularly relevant to settings where online interaction or external supervision is
 17 costly or infeasible, and is conceptually distinct from offline RL approaches that
 18 rely on logged, reward-labeled trajectories.
 19

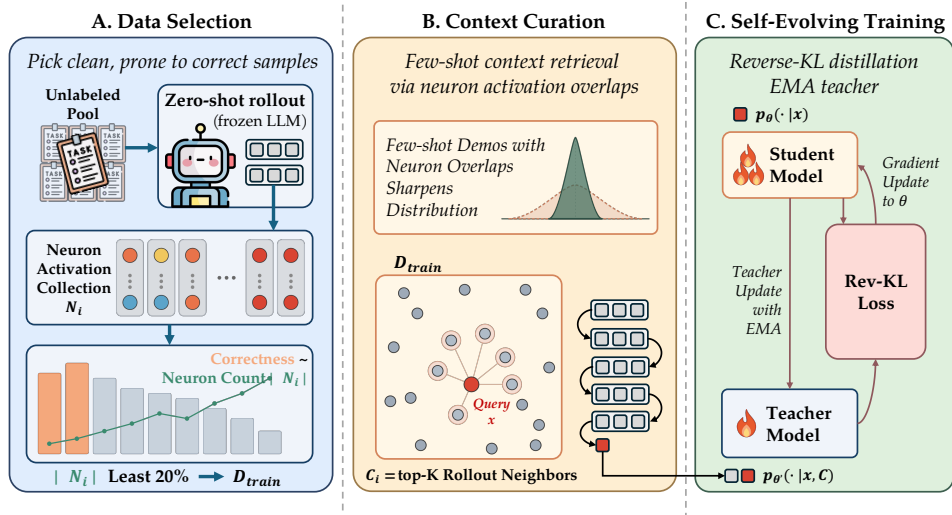


Figure 1: Overview of the Proposed Neuron-OPSD.

20 1 Introduction

21 Recent advances in Large language models (LLMs) have demonstrated remarkable performances
22 in general tasks, however, adapting them to highly specialized domains such as education, law, or
23 STEM subjects, remains difficult, which mainly due to the scarcity of human expert annotations.

24 To mitigate the scarcity of expert annotations, recent work has explored annotation-free self-training,
25 where the model improves itself using only unlabeled data and signals derived from the model
26 itself. SFT-based methods such as LMSI fine-tune the model on its own generated responses or
27 rationales [Huang et al., 2023]. GRPO-based methods instead construct scalar rewards from the
28 model’s own outputs. TTRL extracts pseudo-labels from majority voting over multiple rollouts
29 and optimizes them with the GRPO algorithm [Zuo et al., 2025], while Intuitor uses entropy-based
30 confidence signals as rewards [Zhao et al., 2025b]. Despite their empirical successes, these annotation-
31 free paradigms still face fundamental limitations. SFT-based self-training can suffer from catastrophic
32 forgetting, especially for out-of domain tasks, while GRPO-based optimization faces challenges
33 such as vanishing gradients under low within-group reward variance and entropy collapse during
34 prolonged training. More fundamentally, GRPO collapses each rollout into a scalar reward, yielding
35 only trajectory-level supervision and offering no guidance on which part of the reasoning should be
36 reinforced or revised.

37 On-policy distillation (OPD) offers finer supervision by training the student model against a teacher
38 model’s token-level predictive distribution rather than a hard final-answer reward. Such logit-level
39 targets densify supervision along the generation trajectory and circumvent the learning collapse
40 caused by low within-group reward variances [Song and Zheng, 2026]. However, OPD’s effectiveness
41 relies on the quality of the teacher distribution, which is typically derived from stronger models,
42 verifiers, or curated demonstrations. When such external signals are unavailable, the teacher must be
43 elicited from the model itself through carefully additional information in prompts, i.e contexts around
44 the query. This reframes the problem from designing self-generated rewards alone to deciding which
45 unlabeled samples to train on and which contexts induce a useful teacher–student gap. Crucially, this
46 self-improvement process requires rigorous data selection. Because the model acts as its own teacher,
47 indiscriminately distilling from self-generated signals risks reinforcing its existing miscalibration,
48 spurious reasoning, and hallucinations.

49 While traditional selection methods rely on surface-level metrics like majority-voting or entropy to
50 filter out such unreliable data, these signals often prove insufficient. Recent work has pivoted towards
51 LLMs’ internal dynamics, demonstrating that internal states, particularly neuron activations, exhibit
52 a strong correlation with when a model is genuinely uncertain or hallucinating [Gao et al., 2024,
53 Nostalgebraist, 2020, Chen et al., 2025a, 2026].

54 Motivated by these observations, we propose NEURON-OPSD (N-OPSD), an annotation-free OPD
55 pipeline for LLM self-improvement. In this work, we focus on a two-fold data construction process
56 that determines whether self-distillation provides a useful signal: which unlabeled samples should be
57 used for training, and which few-shot contexts should be used to induce the teacher distribution.

58 As illustrated in Figure 1, NEURON-OPSD utilizes two neuron-based mechanisms. First, NEURON
59 CONSENSUS analyzes the reliability of candidate self-training samples through activation statistics.
60 We measure neuron consensus using the number of activated neurons, finding that this signal is
61 strongly associated with hallucinated or unreliable behaviors of LLMs. However, our results show
62 that consensus alone is not a complete selection criterion. Selecting samples solely by activation
63 count does not consistently yield stronger self-distillation, suggesting that training utility depends on
64 both the reliability of the self-generated signal and the room for teacher-induced sharpening. This
65 suggests that useful training data must be reliable enough to avoid reinforcing hallucinations, yet
66 uncertain enough to leave room for LLMs improvement. Second, NEURON-OPSD employs NEURON
67 OVERLAP for the OPD teacher-context construction. For each query, we retrieve the examples
68 that show the most similar reasoning to the query as demonstrations, which are measured using
69 active neuron representations. Since the teacher is the same model conditioned on an augmented
70 context, these retrieved demonstrations directly shape the token-level distribution used for distillation.
71 Neuron-overlap retrieval, therefore serves as a label-free method to elicit a more informative teacher
72 distribution and create a meaningful teacher-student gap. Together, these two components demonstrate
73 how to effectively leverage internal dynamics for annotation-free OPD. NEURON CONSENSUS serves
74 primarily as a diagnostic signal to identify which samples are reliable or improvable, while NEURON

75 OVERLAP provides a practical mechanism for constructing robust teacher contexts. This makes
76 NEURON-OPSD a practical, data-centric pipeline focused strictly on selecting the optimal data and
77 contexts to form high-quality supervision.

78 We conduct experiments across three datasets covering six source domains. Results show that
79 NEURON-OPSD achieves competitive self-improvement while better preserving cross-domain gener-
80 alization and calibration compared with previous baselines. Further analysis suggests that its gains are
81 driven by both the LLM’s initial capability on the training dataset and the teacher-student distribution
82 gap induced by the contexts. Ablation studies on four SciKnowEval domains further show that neuron
83 consensus data selection is informative but insufficient as a standalone rule, while neuron-overlap
84 retrieval provides an effective context-curation mechanism.

85 Our core contributions are as follows:

- 86 • We utilize NEURON CONSENSUS as a label-free signal for hallucination and data selection. Ac-
87 tivation counts strongly correlate with hallucinated behavior, but selecting neither low- nor high-
88 activation data consistently gives the best self-distillation results. This shows that activation count
89 is informative, but insufficient as the only data-selection rule.
- 90 • We propose NEURON OVERLAP for OPD context curation. It retrieves few-shot demonstrations
91 with similar active-neuron patterns, inducing a more informative teacher distribution and a useful
92 teacher-student gap for annotation-free OPD.
- 93 • We integrate these two signals into NEURON-OPSD. Experiments across specialized-domain
94 benchmarks show that NEURON-OPSD improves self-distillation while better preserving cross-
95 domain performance and reducing calibration collapse.

96 2 Related Work

97 **Annotation-free LLM Self-improvement.** Recent work studies whether LLMs can improve
98 without human labels by deriving supervision from their own outputs or intrinsic signals. Early
99 self-improvement methods bootstrap rationales or self-generated feedback [Zelikman et al., 2022,
100 Yuan et al., 2024, Huang et al., 2023], while recent annotation-free RL methods construct rewards
101 from rollout agreement, entropy, or self-certainty [Zuo et al., 2025, Zhao et al., 2025b]. Other lines
102 explore self-play or zero-data training regimes [Zhao et al., 2025a]. However, recent analysis of
103 unsupervised RLVR shows that intrinsic rewards tend to sharpen the model’s initial distribution and
104 can collapse when confidence is misaligned with correctness [He et al., 2026]. Our work follows this
105 motivation but replaces hard output-level rewards with OPD, using neuron-derived signals to select
106 data and construct teacher contexts without external annotations.

107 **LLM Probing through Internal Neuron Analysis.** Internal-state analysis provides a way to inspect
108 model behavior beyond surface outputs. Early probing methods such as the Logit Lens [Nostalgebraist,
109 2020] and Tuned Lens [Belrose et al., 2023] map intermediate representations into vocabulary space to
110 study how predictions emerge across layers. More recent mechanistic work uses Sparse Autoencoders
111 to disentangle superposed features and obtain more interpretable activation patterns [Gao et al.,
112 2024]. Beyond interpretation, internal activations have also been linked to model reliability. Neuron
113 agreement can signal when models are likely to be correct or hallucinate [Chen et al., 2025b, Cao
114 et al., 2025b], and mechanism-interpretable metrics have been proposed to evaluate utility beyond
115 surface accuracy [Cao et al., 2025a]. Recent work has demonstrated a successful application of
116 neuron-based signals to active few-shot learning, where neuron representations guide the selection
117 of informative examples for annotation, improving few-shot performance with only a small set of
118 high-quality annotated demonstrations Chen et al. [2026]. Our work builds on these findings, using
119 neuron activations not only for analysis but also as label-free signals for data selection and context
120 construction in annotation-free self-improvement.

121 **On-Policy Distillation for LLMs.** OPD trains a student on its own generated trajectories using
122 token-level supervision from a teacher distribution [Ye et al., 2026, Song and Zheng, 2026]. This
123 differs from standard off-policy distillation, where the student learns from fixed teacher-generated
124 data, and from outcome-reward RL, where supervision is usually a scalar reward. Recent work
125 has studied OPD for LLM post-training and self-distillation under settings with stronger teachers,
126 privileged information, or context-conditioned teachers [Hübotter et al., 2026, Zhao et al., 2026, Li
127 et al., 2026]. Our work take OPD as the training objective, and focuses on the annotation-free setting
128 where the teacher context and training data must be constructed without external labels.

129 **3 Preliminaries**

130 To establish the foundation for N-OPSD, we formalize the annotation-free self-evolution setting and
 131 the mechanics of OPD [Ye et al., 2026, Hübötter et al., 2026].

132 **3.1 Problem Formulation**

133 We consider an annotation-free self-evolution scenario. Given a base LLM with parameters θ and
 134 a domain-specific unlabeled dataset $\mathcal{D} = \{x_i\}_{i=1}^N$ of input queries, the objective is to improve the
 135 base model’s performance on the target domain without access to ground-truth labels or any external
 136 model, i.e a different stronger teacher model. We frame this as annotation-free post-training from
 137 a fixed unlabeled prompt pool: the only source of new signal during training is the model’s own
 138 on-policy rollouts and intrinsic activation patterns. No environment interaction, human feedback,
 139 or external oracle is available. This is distinct from offline RL setups, which learn from logged
 140 reward-labeled trajectories, and from online RL post-training such as RLHF, which assumes access
 141 to live preference or reward queries.

142 **3.2 On-Policy Distillation**

143 On-Policy Distillation, [Song and Zheng, 2026] trains the base or student model π_θ by minimizing a
 144 per-token reverse KL loss to a teacher model π_{teacher} on trajectories sampled from the student model:

$$\mathcal{L}_{\text{OPD}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot | x)} \left[\sum_{t=1}^{|y|} \text{KL}(\pi_\theta(\cdot | x, y_{<t}) \parallel \pi_{\text{teacher}}(\cdot | x, y_{<t})) \right], \quad (1)$$

145 The effectiveness of the student model largely depends on the quality of π_{teacher} , which is typically a
 146 larger model or the same model augmented with additional context, e.g., conditioning on external
 147 feedback like compiler traces [Hübötter et al., 2026]. In our annotation-free setting, both teacher
 148 and student model are initialized from the same base LLM and no external feedback is available,
 149 so constructing a reliable and high-quality π_{teacher} becomes the key challenge. To address this, we
 150 condition the teacher model on auxiliary context c constructed from unlabeled domain data alone,
 151 without any additional information, improving its next-token distribution on the target task, i.e.,
 152 $\pi_{\text{teacher}}(\cdot | x, y_{<t}) := \pi_\theta(\cdot | x, c, y_{<t})$.

153 **3.3 Internal Signals for Data Selection**

154 The additional context c for the teacher model is constructed primarily through few-shot examples
 155 with the model’s own generated solutions. Rather than using external supervision to verify the
 156 correctness of these pseudolabels, we probe the model’s internal dynamics as a proxy for correctness.
 157 Now we introduce the calculation for the active neuron set.

158 **Internal Neuron Signal Extraction.** Building on the logit-lens view that intermediate representa-
 159 tions can be read through the LM head [Nostalgebraist, 2020], and recent evidence that neuron-level
 160 activation features reflect response reliability [Cao et al., 2025b, Chen et al., 2025a, 2026], we
 161 follow Chen et al. [2026] to extract the activated MLP neurons detailed as follows. For layer l , let
 162 \mathbf{h}^l be the MLP input, \mathbf{W}_{in}^l and $\mathbf{W}_{\text{out}}^l$ be the up- and down-projection matrices. We compute the
 163 activation as:

$$\mathbf{k}^l = \sigma(\mathbf{h}^l \mathbf{W}_{\text{in}}^l), \quad (2)$$

164 Next, given a generated token \hat{y} , we score each neuron at layer l and index i by projecting its
 165 downstream contribution to the unembedding space:

$$S_{\hat{y},i}^l = k_i^l \cdot (\mathbf{w}_{\text{out},i}^l \cdot \mathbf{e}_{\hat{y}}), \quad (3)$$

166 where $\mathbf{w}_{\text{out},i}^l$ is the i -th row of $\mathbf{W}_{\text{out}}^l$ and $\mathbf{e}_{\hat{y}}$ is the unembedding vector for \hat{y} . This early-unembedding
 167 score filters for neurons that actively promote the model’s generated token. For each generated token
 168 y_i , we define the activated neurons by keeping the top- K scored neurons across all layers. We then
 169 aggregate the retained neurons across the generated tokens in the response by taking their union
 170 set. The resulting sparse set, denoted as $\mathcal{N}(x)$, serves as the internal signal used later for consensus
 171 estimation and context retrieval.

¹We follow Chen et al. [2026] to set K as 5000.

172 4 Methodology

173 We present a neuron-activation-guided framework for selecting training data for self-distillation and
174 curate contexts for self-evolving LLMs. Given an unlabeled data pool $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$, our
175 approach selects a subset $\mathcal{S} \subset \mathcal{D}$ for self-distillation training, without requiring ground-truth labels.
176 The framework consists of four stages: neuron activation extraction, sample selection based on
177 neuron-consensus, context curation via Jaccard-nearest few-shot retrieval, and self-evolving training.

178 4.1 Neuron Activation Extraction

179 We instantiate the internal signal defined in §3.3 for every sample in the unlabeled pool. Concretely,
180 we collect rollouts with the model in zero-shot, compute early-unembedding contribution scores for
181 MLP neurons, and store the set of activated neurons $\mathcal{N}(x)$ as the sample’s sparse activation pattern.

182 4.2 Sample Selection via Neuron Consensus

183 Given the activation pattern $\mathcal{N}(x)$, we measure neuron consensus of x by the size of its activation set:

$$s(x) = |\mathcal{N}(x)|, \quad (4)$$

184 where a smaller $s(x)$ indicates higher consensus, and a larger $s(x)$ indicates diffuse activation that
185 could related to hallucinations Chen et al. [2025b, 2026]. We interpret $s(x)$ as an annotation-free
186 proxy for the model’s sample-level reliability. Since samples with lower $s(x)$ tend to be answered
187 more correctly, applying OPD-based self-improvement to these samples helps reinforce the model’s
188 already-correct reasoning and predictions, helping to avoid exacerbating hallucinations or self-
189 confirmation bias.

190 4.3 Neuron Overlap for Teacher Context Construction

191 To construct relevant few-shot contexts for the self-distillation teacher, we measure the overlap
192 between neuron activation patterns using the Jaccard distance:

$$J(x_i, x_j) = 1 - \frac{|\mathcal{N}(x_i) \cap \mathcal{N}(x_j)|}{|\mathcal{N}(x_i) \cup \mathcal{N}(x_j)|}, \quad (5)$$

193 A small Jaccard distance indicates that the model processes two samples using similar neurons,
194 suggesting they share similar knowledge circuits or go over similar reasoning path. For each query
195 sample x_q , we construct its K -shot demonstration set by selecting the K nearest neighbors by:

$$\mathcal{C}_K(x_q) = \arg \min_{\substack{S \subset \mathcal{D} \setminus \{x_q\} \\ |S|=K}} \sum_{x_j \in S} J(x_q, x_j), \quad (6)$$

196 This selection ensures that few-shot examples are processed through similar neural pathways as the
197 query, providing the teacher with the most relevant contextual knowledge for distillation, helping
198 sharpening its prediction distribution.

199 4.4 Self-Improvement via On-Policy Distillation

200 To enable annotation-free self-improvement, we integrate neuron-guided sample selection and context
201 curation with OPD. The student is the zero-shot policy parameterized by the current model θ , while
202 the teacher is an EMA-updated copy $\bar{\theta}$ conditioned on the additional context $\mathcal{C}_K(x)$:

$$\pi_{\bar{\theta}}^{\text{stu}}(\cdot | x, y_{<t}) := \pi_{\theta}(\cdot | x, y_{<t}), \quad (7)$$

$$q_{\bar{\theta}}^{\text{tea}}(\cdot | \mathcal{C}_K(x), x, y_{<t}) := \pi_{\bar{\theta}}(\cdot | \mathcal{C}_K(x), x, y_{<t}), \quad (8)$$

$$\bar{\theta} \leftarrow (1 - \tau)\bar{\theta} + \tau\theta, \quad (9)$$

203 where $\mathcal{C}_K(x)$ denotes the additional context guided by neuron-overlap retrieved for sample x . During
204 optimization, the teacher distribution is treated as a fixed target, and gradients are applied only to the
205 student policy $\pi_{\bar{\theta}}^{\text{stu}}$. The teacher is updated by exponential moving average (EMA) rate τ .

206 The OPD training objective minimizes the reverse KL divergence on student-generated rollouts:

$$\mathcal{L}_{\text{OPD}} = \mathbb{E}_{x \sim \mathcal{S}, y_{<t} \sim \pi_{\bar{\theta}}^{\text{stu}}} \left[\sum_{t=1}^T D_{\text{KL}}(\pi_{\bar{\theta}}^{\text{stu}}(\cdot | x, y_{<t}) \| q_{\bar{\theta}}^{\text{tea}}(\cdot | c_x, x, y_{<t})) \right]. \quad (10)$$

207 where \mathcal{S} denotes the selected subset, and $c_x = \mathcal{C}_K(x)$ is the neuron-overlap context retrieved for
 208 sample x . During optimization, the teacher distribution is treated as a fixed target, and gradients
 209 are applied only to the student policy. The teacher parameters $\bar{\theta}$ are updated by exponential moving
 210 average of the student parameters, yielding a temporally smoothed model updates.

211 This objective transfers the distributional benefits of neuron-aligned few-shot contexts into the
 212 zero-shot student policy, internalizing context-induced distributions through parameter updates and
 213 eliminating reliance on demonstrations at inference. Crucially, the entire pipeline is annotation-free:
 214 neuron activations guide both sample selection and context curation, while N-OPSD uses the model’s
 215 own context-conditioned predictions as supervision. Algorithm 1 summarizes the full procedure.

216 5 Experiments

217 To establish the impact of each proposed component, we first analyze the effectiveness of each
 218 individually, namely, neuron consensus, neuron overlap, and the OPD training objective when
 219 combined with neuron consensus and neuron overlap. For the analyses in Section 5.1 and Section 5.2,
 220 we use **SciKnowEval** [Feng et al., 2025], a multiple-choice benchmark spanning four scientific
 221 domains, *Biology*, *Material*, *Physics*, and *Chemistry*, each split into 80% train and 20% test sets.

222 5.1 Analysis: Neuron Consensus based Data Selection

223 Does Neuron Consensus Correlate with Performance?

224 We validate the relationship between the activation density, i.e., the total count of activated neurons $|\mathcal{N}(x)|$, and the
 225 model’s hallucination rates. As illustrated in Figure 2, our analysis across all four SciKnowEval domains reveals a
 226 strong, robust monotonic correlation: queries that activate a larger number of neurons are significantly more likely to
 227 result in incorrect rollouts.
 228
 229
 230

231 This phenomenon is attributed to Neuron Consensus. Sparser activations with lower $|\mathcal{N}(x)|$ indicate a higher
 232 degree of consensus within the network. The model effortlessly retrieves a clear, unconflicted knowledge trace. Con-
 233 versely, a denser, widespread activation pattern suggests internal conflict and uncertainty, where the model strug-
 234 gles to reconcile disparate neural pathways. This finding fundamentally addresses our data selection bottleneck: by
 235 isolating queries with the lowest activation counts, we can filter out hallucination-prone samples and secure pristine
 236 targets for self-supervision, even without external labels.
 237
 238
 239
 240
 241

242 **Are Easy Samples Better Training Data for OPD Training?** We rank the source pool by neuron
 243 count $s(x)$ and train on the subset that ranked top and bottom 20% following the OPD framework.
 244 For evaluation, we sample 8 rollouts for each test query and report two metrics: **Avg@8**, the mean
 245 per-sample accuracy; and **ECE**, the expected calibration error over the per-query majority-vote
 246 confidence with 15 equal-width bins, lower the better. Table. 1 reports the results on SciKnowEval.

247 Intuitively, subsets with higher activation numbers, i.e. Top-20%, will generally noisier than the
 248 subsets with lower activation numbers, i.e., Bottom-20%. As lower activation numbers indicates
 249 a lower hallucination rates. However, in training ablations, with OPD, Top-20% can still have
 250 reasonable improvements on some domains such as Bio. and Mat., also with a competitive calibration.

251 Furthermore, a trade-off exists between learning utility and noise levels. Specifically, samples prone
 252 to hallucination theoretically offer higher learning value since they represent areas where the LLM
 253 consistently fails. However, in self-improvement scenarios, the model struggles to generate reliable

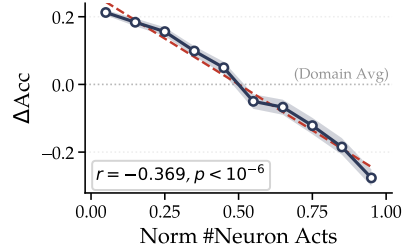


Figure 2: Neuron consensus correlates accuracy. ΔAcc is the gap between average domain accuracy and bin accuracy, while $\#\text{Neuron Acts}$ is normalized domain-wise. Separate domain-wise relational results shown in Figure 3.

Selection	Avg@8 (\uparrow)				ECE (\downarrow)			
	CHEM.	BIO.	MAT.	PHYS.	CHEM.	BIO.	MAT.	PHYS.
Qwen3-4B	72.11	73.98	71.12	80.20	0.173	0.195	0.213	0.095
+Top-20%	68.46 _{-3.65}	74.59 _{+0.61}	73.82 _{+2.70}	82.00 _{+1.80}	0.133 _{-0.040}	0.198 _{+0.003}	0.194 _{-0.019}	0.116 _{+0.021}
+Bottom-20%	72.03 _{-0.08}	74.02 _{+0.04}	73.39 _{+2.27}	83.13 _{+2.93}	0.171 _{-0.002}	0.193 _{-0.002}	0.207 _{-0.006}	0.114 _{+0.019}

Table 1: Analysis results of consensus-based selection, reported in-domain Avg@8 and ECE.

254 supervision signals for these difficult cases. Conversely, while the LLM can easily produce correct self-
255 generated signals for simpler samples, their contribution to overall improvement is marginal. Given
256 this persistent trade-off, training on a refined top-20% subset can still yield substantial performance
257 gains, even in the presence of noise.

258 5.2 Analysis: Neuron Overlap guided Context Curation

259 In this experiment, we fix the Bottom-20% as OPD training data, we compare neuron-overlap retrieval
260 with uniform in-domain random retrieval. As shown in Table. 2, NEURON-JACCARD improves over
261 the base model on three of the four SciKnowEval sources and outperforms random retrieval on CHEM.
262 and PHYS.. Random in-domain retrieval is also competitive on BIO. and MAT., indicating that
263 in-domain demonstrations already provide a strong context prior.

264 We further analyze why the marginal gain over random varies across domains. Because both random
265 and neuron-based retrieval are restricted to the same source domain, the difference depends on
266 how discriminative the within-domain contexts are. We measure reasoning diversity by the mean
267 pairwise neuron-Jaccard distance between training samples: CHEM. 0.870, BIO. 0.634, MAT. 0.628,
268 and PHYS. 0.640. In more diverse domains, such as CHEM., neuron-overlap retrieval can select
269 demonstrations that match the query-specific reasoning mode, yielding a clear advantage over random
270 selection. In more homogeneous domains, many examples activate similar neuron sets, so nearest-
271 neighbor and random in-domain contexts become partially redundant, reducing the marginal benefit
272 of Jaccard retrieval. This explains why random retrieval remains competitive on BIO. and MAT.,
273 while also showing that NEURON OVERLAP is most useful when neuron similarity can distinguish
274 meaningful reasoning modes within the source pool.

Retriever	CHEM.	BIO.	MAT.	PHYS.
Qwen3-4B-Thinking	72.11	73.98	71.12	80.20
+Random	68.48 _{-3.63}	74.14 _{+0.16}	73.79 _{+2.67}	81.97 _{+1.77}
+Neuron-Jaccard (ours)	72.03 _{-0.08}	74.02 _{+0.04}	73.39 _{+2.27}	83.13 _{+2.93}

Table 2: Analysis results of the few-shot retriever, reported in-domain Avg@8 on SciKnowEval.

275 5.3 Experimental Setup

276 We then perform the evaluations on various self-improvement training algorithms. As we didn’t
277 observe a significant performance and calibration difference of using the subsets with the bottom 20%
278 vs the top 20% #neuron activations, therefore we demonstrated the bottom 20% results.

279 **Datasets.** We evaluate on three different datasets. **SciKnowEval** [Feng et al., 2025] contains
280 four scientific multiple-choice domains: *Biology*, *Material*, *Physics*, and *Chemistry*, each split
281 80/20 for training and testing. **Edu-Feedback** [Wu and Schunn, 2023] is a binary feedback-quality
282 classification dataset with 1,799 training and 1,000 testing samples. **MMLU-Pro** [Wang et al., 2024]
283 is a challenging multi-domain multiple-choice benchmark. We split 80% data for self-improvement.

284 **Baselines and Models.** We compare with both the off-policy and on-policy self-improvement
285 baselines. Including SFT-based **LMSI** [Huang et al., 2023], on-policy methods **TTRL** [Zuo et al.,
286 2025], and **Intuitor** [Zhao et al., 2025b], an annotation-free RL method using self-certainty as reward.
287 We evaluate on *Qwen3-4B-Thinking-2507* [Yang et al., 2025].

288 **Metrics.** We sample 8 rollouts for each test query and report three metrics: **Avg@8**, the mean
289 per-sample accuracy; **Maj@8**, the majority-vote accuracy; and **ECE**, the expected calibration error
290 over the per-query majority-vote confidence with 15 equal-width bins, lower the better.

291 **5.4 Performance Results**

292 To more comprehensively evaluate the efficacy of the self-improving baselines and the proposed
 293 method, we focus on both in-domain and cross-domain performance gain and calibration variances.
 294 Specifically, we train each method on a single source domain and evaluate on all target test sets.
 295 Table. 3 and Table. 4 report Accuracy and ECE on Qwen3-4B-Thinking-2507.

Method	Eval.	BIO	MAT.	SciKnowEval			Avg.	Edu.	MMLU-Pro	Avg.
				PHYS.	CHEM					
Qwen3-4B base		73.98	71.12	80.20	72.11	74.35	64.95	72.16	72.42	
LMSI	in	74.59 ^{+0.61}	69.08 ^{-2.04}	78.54 ^{-1.66}	73.72 ^{+1.61}	73.98 ^{-0.37}	66.95 ^{+2.00}	68.60 ^{-3.56}	71.91 ^{-0.51}	
	cross	68.57 ^{-3.54}	65.48 ^{-7.20}	67.02 ^{-3.84}	71.19 ^{-1.29}	68.06 ^{-3.97}	58.14 ^{-15.77}	73.46 ^{+0.99}	67.31 ^{-5.11}	
TTRL	in	75.00 ^{+1.02}	73.26 ^{+2.14}	82.98 ^{+2.78}	72.43 ^{+0.32}	75.92 ^{+1.56}	64.95 ^{+0.00}	72.04 ^{-0.12}	73.44 ^{+1.02}	
	cross	73.31 ^{+1.20}	73.32 ^{+0.64}	71.28 ^{+0.41}	73.46 ^{+0.97}	72.84 ^{+0.81}	74.47 ^{+0.55}	73.76 ^{+1.29}	73.27 ^{+0.84}	
Intuitor ²	in	76.21 ^{+2.23}	73.41 ^{+2.29}	83.43 ^{+3.23}	72.63 ^{+0.52}	76.42 ^{+2.07}	—	71.91 ^{-0.25}	75.52 ^{+1.60}	
	cross	69.50 ^{-2.61}	72.79 ^{+0.11}	68.33 ^{-2.53}	70.76 ^{-1.72}	70.35 ^{-1.69}	—	73.30 ^{+0.83}	70.94 ^{-1.18}	
N-OPSD	in	74.02 ^{+0.04}	73.39 ^{+2.27}	83.13 ^{+2.93}	72.03 ^{-0.08}	75.64 ^{+1.29}	72.19 ^{+7.24}	72.04 ^{-0.12}	74.47 ^{+2.05}	
	cross	72.71 ^{+0.61}	73.27 ^{+0.59}	71.37 ^{+0.51}	73.59 ^{+1.11}	72.74 ^{+0.70}	75.22 ^{+1.31}	73.74 ^{+1.27}	73.32 ^{+0.90}	

Table 3: Cross-domain evaluation on Qwen3-4B-Thinking-2507, Avg@8. For each method, *in* reports the score of the model trained and evaluated on the same source domain, while *cross* reports the average score of the same source-trained model on the other target domains. Subscripts report Δ against the untrained Qwen3-4B baseline computed on the same target-domain support. The full source-target results are provided in Table 7.

Method	Eval.	BIO	MAT.	SciKnowEval			Avg.	Edu.	MMLU-Pro	Avg.
				PHYS.	CHEM					
Qwen3-4B base		0.195	0.213	0.095	0.173	0.169	0.246	0.184	0.184	
LMSI	in	0.202 ^{+0.007}	0.225 ^{+0.012}	0.111 ^{+0.016}	0.179 ^{+0.006}	0.179 ^{+0.010}	0.255 ^{+0.009}	0.171 ^{-0.013}	0.191 ^{+0.007}	
	cross	0.171 ^{-0.011}	0.179 ^{+0.001}	0.202 ^{+0.000}	0.182 ^{-0.004}	0.184 ^{-0.004}	0.095 ^{-0.077}	0.162 ^{-0.022}	0.165 ^{-0.019}	
TTRL	in	0.204 ^{+0.009}	0.208 ^{-0.005}	0.123 ^{+0.028}	0.192 ^{+0.019}	0.182 ^{+0.013}	0.248 ^{+0.002}	0.198 ^{+0.014}	0.196 ^{+0.012}	
	cross	0.186 ^{+0.004}	0.187 ^{+0.009}	0.204 ^{+0.002}	0.199 ^{+0.012}	0.194 ^{+0.007}	0.178 ^{+0.006}	0.190 ^{+0.006}	0.191 ^{+0.007}	
Intuitor	in	0.208 ^{+0.013}	0.212 ^{-0.001}	0.117 ^{+0.022}	0.198 ^{+0.026}	0.184 ^{+0.015}	—	0.174 ^{-0.010}	0.182 ^{+0.010}	
	cross	0.226 ^{+0.044}	0.187 ^{+0.008}	0.234 ^{+0.032}	0.224 ^{+0.038}	0.218 ^{+0.030}	—	0.185 ^{+0.000}	0.211 ^{+0.024}	
N-OPSD	in	0.193 ^{-0.003}	0.207 ^{-0.006}	0.114 ^{+0.019}	0.171 ^{-0.002}	0.171 ^{+0.002}	0.191 ^{-0.056}	0.180 ^{-0.004}	0.176 ^{-0.008}	
	cross	0.177 ^{-0.005}	0.177 ^{-0.001}	0.194 ^{-0.008}	0.187 ^{+0.000}	0.184 ^{-0.004}	0.177 ^{+0.005}	0.186 ^{+0.002}	0.183 ^{-0.001}	

Table 4: Cross-domain calibration on Qwen3-4B-Thinking-2507, reported as ECE, where lower is better, computed from per-query majority-vote confidence. For each method, *in* reports the ECE of the model trained and evaluated on the same source domain, while *cross* reports the average ECE of the same source-trained model on the other target domains. Subscripts report Δ against the untrained Qwen3-4B baseline computed on the same target-domain support, where negative values indicate improved calibration. The full results are provided in Table 8.

296 **Performance.** N-OPSD improves source-domain performance on four of the six training sources
 297 in Table. 3, with especially clear gains on MAT., PHYS., and EDU. On MAT., PHYS., and EDU.,
 298 N-OPSD matches or exceeds TTRL’s source-domain gain, and its averaged row attains the highest
 299 overall Avg@8 among the compared methods. The main advantage is not uniform source-domain
 300 dominance, but a stronger accuracy-preservation trade-off: every SciKnow-trained N-OPSD model
 301 improves EDU. over the base model, whereas TTRL is mixed. LMSI and Intuitor often suffer
 302 large drops. Thus, N-OPSD delivers source-domain improvement where the teacher provides useful
 303 sharpening while better preserving non-source capability.

304 **Calibration.** Table. 4 shows that TTRL and Intuitor generally increase average ECE, while N-OPSD
 305 keeps the calibration cost smaller and slightly reduces overall ECE in the averaged row. The strongest
 306 calibration effect appears on EDU.: every N-OPSD row lowers EDU. ECE relative to the base
 307 model, whereas TTRL is mixed and Intuitor often worsens it substantially. Several N-OPSD rows
 308 simultaneously improve overall Avg@8 and reduce overall ECE, and N-OPSD_{EDU.} gives the largest
 309 overall accuracy gain together with the lowest overall ECE among N-OPSD rows. These results show
 310 where the method works best, it improves accuracy without the broad calibration inflation seen in
 311 reward-based annotation-free RL, although it can still increase ECE on some SciKnowEval targets.

²The Intuitor training on the Edu. dataset yields no improvement across validation steps, thus we drop the results.

312 **5.5 What drives the self-improvement, and when does it benefit most?**

313 In terms of self-training, He et al. [2026] characterise intrinsic self-supervision methods as sharpening
 314 the model’s prior distribution rather than adding new information, which succeeds only when that
 315 prior is already aligned with correctness. Our consensus stratification makes the trade concrete at
 316 the sample level, with high-consensus samples supplying a reliable prior to sharpen against and
 317 low-consensus samples supplying an unreliable one. Meanwhile, OPD derives its training signal
 318 from a soft teacher-student distributional gap rather than from hard scalar rewards. This allows part
 319 of the teacher’s uncertainty to be preserved during distillation, yielding a more favorable trade-off
 320 between calibration and correctness. Here, we conduct a data-centric analysis of the selected training
 321 pool to examine how data dynamics drive the gains of OPD-based LLM self-training.

Domain	Base Test Avg@8	Train Maj@8	Train Avg@8	Maj–Avg	in-domain gain
CHEM.	72.11	98.23	98.18	0.04	−0.08
BIO.	73.98	98.73	98.64	0.09	+0.04
MAT.	71.12	91.93	91.59	0.34	+2.27
PHYS.	80.20	96.31	96.27	0.04	+2.93

Table 5: Base rollout statistics on the training pool versus N-OPSD’s in-domain Avg@8 gain. All statistics are computed with Qwen3-4B-Thinking-2507: Test Avg@8 is measured on the held-out test set, while training-pool statistics are measured on the bottom-20% subset by neuron count.

Domain	H_{zs} Student	H_{neur} Teacher	ΔH T–S	Avg@8 Gain
CHEM.	0.240	0.237	−0.002	−0.08
BIO.	0.238	0.240	+0.001	+0.04
MAT.	0.268	0.251	−0.017	+2.27
PHYS.	0.257	0.247	−0.010	+2.93

Table 6: Teacher-student per-token logprob entropy on the selected training pool. The teacher conditions on $K=10$ neuron-Jaccard nearest demonstrations, and the student is zero-shot. $\Delta H < 0$ means the teacher is sharper at the token level. The two domains with significant teacher sharpening, MAT. and PHYS., are also the two with the largest in-domain Avg@8 gain.

322 **Sharpening Room on High-Consensus Samples.** A necessary precondition for N-OPSD to deliver
 323 a meaningful in-domain gain is that the high neuron-consensus pool used for distillation still contains
 324 enough gap for sharpening, i.e., carries disagreement across the rollouts. Table. 5 reports the $n=8$
 325 rollout stats of the untrained model on this pool. Across CHEM., BIO., and MAT., the Maj–Avg
 326 gap on the high-consensus training pool tracks N-OPSD’s in-domain gain: larger residual rollout
 327 disagreement gives OPD distillation more uncertainty to resolve. However, PHYS. is already near-
 328 unanimous and still obtains a large gain, showing that vote-level sharpening room alone is not
 329 sufficient to explain all improvements.

330 **Teacher-Student Gap Drives Gain.** The sharpening-room analysis captures part of when N-OPSD
 331 has signal to exploit, but not whether the few-shot teacher provides a useful target. On the selected
 332 pool, vote-level metrics saturate, as both the zero-shot student and the neuron-Jaccard teacher reach
 333 $\sim 98\%$ Avg@8 in Table. 5, so the teacher-student gap must be measured at the token level. Table. 6
 334 reports per-token logprob entropy, averaged over the response trajectory and paired across queries.
 335 On MAT. and PHYS. the teacher meaningfully sharpens the token distribution, with $\Delta H = -0.017$
 336 and -0.010 , and these are exactly the two domains where N-OPSD delivers the strongest in-domain
 337 gains of $+2.27$ and $+2.93$. On BIO. and CHEM. the teacher provides no token-level sharpening, with
 338 $|\Delta H| < 0.003$ and no significant difference, and the in-domain gain collapses correspondingly.

339 **6 Conclusion**

340 In this work, we propose NEURON-OPSD for annotation-free post-training from a fixed unlabeled
 341 prompt pool, where no ground-truth annotations are available. The method uses Neuron Consensus to
 342 select more reliable self-improvement data and Neuron Overlap to construct neuron-aligned teacher
 343 contexts for on-policy distillation. These two components are proposed for the central challenge of
 344 unreliable self-supervision, which can otherwise amplify hallucinations, miscalibration, or negative
 345 transfer. Across our evaluations, NEURON-OPSD achieves competitive in-domain gains while better
 346 preserving cross-domain capability and calibration than prior annotation-free methods.

347 Limitations

348 There are several limitations remaining in this work. First, in-domain gains are not uniform across
349 domains: BIO. and CHEM. show negligible improvement, consistent with Tab. 6 showing no
350 teacher token-distribution sharpening on those domains. Second, although NEURON CONSENSUS
351 correlates with hallucinated behavior, our top- and bottom-activation ablations show that activation
352 count alone is insufficient for data selection. This poses a fundamental question: how do we
353 balance self-improvement against signal reliability? Identifying the precise conditions under which
354 LLMs can be trusted to self-supervise is left for future work. Third, NEURON OVERLAP’s retrieval
355 coherence depends on domain reasoning diversity; as the case study in App. C.2 shows, domains
356 with homogeneous reasoning patterns can make neuron-Jaccard retrieval less discriminative and
357 push contexts toward generic examples. Fourth, although we compare against uniform in-domain
358 random few-shot retrieval, we do not compare against other alternative teacher-context construction
359 strategies, such as retrieval-augmented or prompt-engineered contexts, that might produce similar
360 teacher sharpening without neuron overlap. Also, the evaluation covers three datasets, SciKnowEval,
361 Edu-Feedback, and MMLU-Pro, so the observed side-effect tradeoffs may change under broader
362 domain coverage and different context-construction mechanisms.

363 References

- 364 Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella
365 Biderman, and Jacob Steinhardt. Eliciting latent predictions from transformers with the tuned lens.
366 *arXiv preprint arXiv:2303.08112*, 2023.
- 367 Yixin Cao, Jiahao Ying, Yaoning Wang, Xipeng Qiu, Xuanjing Huang, and Yugang Jiang. Model
368 utility law: Evaluating llms beyond performance through mechanism interpretable metric, 2025a.
369 URL <https://arxiv.org/abs/2504.07440>.
- 370 Yixin Cao, Jiahao Ying, Yaoning Wang, Xipeng Qiu, Xuanjing Huang, and Yugang Jiang. Model
371 utility law: Evaluating LLMs beyond performance through mechanism interpretable metric, 2025b.
372 URL <http://arxiv.org/abs/2504.07440>.
- 373 Kang Chen, Yaoning Wang, Kai Xiong, Zhuoka Feng, Wenhe Sun, Haotian Chen, and Yixin Cao. Do
374 llms signal when they’re right? evidence from neuron agreement. *arXiv preprint arXiv:2510.26277*,
375 2025a.
- 376 Kang Chen, Yaoning Wang, Kai Xiong, Zhuoka Feng, Wenhe Sun, Haotian Chen, and Yixin
377 Cao. Do LLMs signal when they’re right? evidence from neuron agreement, 2025b. URL
378 <http://arxiv.org/abs/2510.26277>.
- 379 Zhuowei Chen, Liwei Chen, Christian Schunn, Raquel Coelho, and Xiang Lorraine Li. Neuron-aware
380 active few-shot learning for LLMs. In *Proceedings of the 64th Annual Meeting of the Association
381 for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics,
382 2026.
- 383 Kehua Feng, Xinyi Shen, Weijie Wang, Xiang Zhuang, Yuqi Tang, Qiang Zhang, and Keyan Ding.
384 Sciknoweval: Evaluating multi-level scientific knowledge of large language models, 2025. URL
385 <https://arxiv.org/abs/2406.09098>.
- 386 Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya
387 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint
388 arXiv:2406.04093*, 2024.
- 389 Bingxiang He, Yuxin Zuo, Zeyuan Liu, Shangziqu Zhao, Zixuan Fu, Junlin Yang, Cheng Qian, Kaiyan
390 Zhang, Yuchen Fan, Ganqu Cui, Xiusi Chen, Youbang Sun, Xingtai Lv, Xuekai Zhu, Li Sheng,
391 Ran Li, Huan-ang Gao, Yuchen Zhang, Bowen Zhou, Zhiyuan Liu, and Ning Ding. How far can
392 unsupervised RLVR scale LLM training?, 2026. URL <http://arxiv.org/abs/2603.08660>.
- 393 Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han.
394 Large language models can self-improve. In Houda Bouamor, Juan Pino, and Kalika Bali, editors,
395 *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages
396 1051–1068, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/
397 v1/2023.emnlp-main.67. URL <https://aclanthology.org/2023.emnlp-main.67/>.

- 398 Jonas Hübötter, Frederike Lübeck, Lejs Behric, Anton Baumann, Marco Bagatella, Daniel Marta,
399 Ido Hakimi, Idan Shenfeld, Thomas Kleine Buening, Carlos Guestrin, and Andreas Krause.
400 Reinforcement learning via self-distillation, 2026. URL <http://arxiv.org/abs/2601.20802>.
- 401 Yaxuan Li, Yuxin Zuo, Bingxiang He, Jinqian Zhang, Chaojun Xiao, Cheng Qian, Tianyu Yu, Huan-
402 ang Gao, Wenkai Yang, Zhiyuan Liu, et al. Rethinking on-policy distillation of large language
403 models: Phenomenology, mechanism, and recipe. *arXiv preprint arXiv:2604.13016*, 2026.
- 404 Nostalgebraist. Interpreting GPT: The logit lens. [https://www.lesswrong.com/posts/
405 AcCRPn5w6xmtPRD7j/interpreting-gpt-the-logit-lens](https://www.lesswrong.com/posts/AcCRPn5w6xmtPRD7j/interpreting-gpt-the-logit-lens), 2020. LessWrong.
- 406 Mingyang Song and Mao Zheng. A survey of on-policy distillation for large language models, 2026.
407 URL <https://arxiv.org/abs/2604.00626>.
- 408 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
409 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging
410 multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024.
- 411 Yong Wu and Christian D Schunn. Passive, active, and constructive engagement with peer feedback:
412 A revised model of learning from peer feedback. *Contemporary Educational Psychology*, 73:
413 102160, 2023.
- 414 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
415 Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*,
416 2025.
- 417 Tianzhu Ye, Li Dong, Xun Wu, Shaohan Huang, and Furu Wei. On-policy context distillation for
418 language models, 2026. URL <http://arxiv.org/abs/2602.12275>.
- 419 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and
420 Jason Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2024.
- 421 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. STaR: Bootstrapping reasoning with
422 reasoning, 2022. URL <http://arxiv.org/abs/2203.14465>.
- 423 Andrew Zhao, Yiran Wu, Yang Yue, Tong Wu, Quentin Xu, Matthieu Lin, Shenzhi Wang, Qingyun
424 Wu, Zilong Zheng, and Gao Huang. Absolute zero: Reinforced self-play reasoning with zero data.
425 *arXiv preprint arXiv:2505.03335*, 2025a.
- 426 Siyan Zhao, Zhihui Xie, Mengchen Liu, Jing Huang, Guan Pang, Feiyu Chen, and Aditya Grover.
427 Self-distilled reasoner: On-policy self-distillation for large language models, 2026. URL [http://
428 arxiv.org/abs/2601.18734](http://arxiv.org/abs/2601.18734).
- 429 Xuandong Zhao, Zhewei Kang, Aosong Feng, Sergey Levine, and Dawn Song. Learning to reason
430 without external rewards. *arXiv preprint arXiv:2505.19590*, 2025b.
- 431 Yuxin Zuo, Kaiyan Zhang, Li Sheng, Shang Qu, Ganqu Cui, Xuekai Zhu, Haozhan Li, Yuchen Zhang,
432 Xinwei Long, Ermo Hua, Biqing Qi, Youbang Sun, Zhiyuan Ma, Lifan Yuan, Ning Ding, and
433 Bowen Zhou. TTRL: Test-time reinforcement learning, 2025. URL [http://arxiv.org/abs/
434 2504.16084](http://arxiv.org/abs/2504.16084).

435 A Implementation Details

436 **Neuron contribution computation.** We register forward hooks on the activation function of each
437 transformer MLP layer, `model.layers[1].mlp.act_fn`. For each response token position t , we
438 capture the post-activation hidden state $\mathbf{a}_{i,t} \in \mathbb{R}^{d_{\text{inner}}}$ and compute contribution scores via Eq. 3.
439 We retain the top 2,000 neurons per layer per chunk, then apply global Top- K deduplication with
440 $K = 5,000$ across all layers. When the same layer-neuron pair appears across multiple chunks,
441 namely groups of response positions, we keep the maximum contribution score.

442 **N-OPSD training configuration.** All N-OPSD models are trained using the veRL framework with
443 Ray-based FSDP distributed training across 4 GPUs, batch size 16, micro-batch size 4, and 150
444 training steps. We use reverse KL divergence and EMA teacher updates with rate 0.01.

445 **TTRL training configuration.** TTRL models follow the standard open-source setup of Zuo et al.
 446 [2025], using GRPO with majority-vote pseudo-labels. We train with 4 GPUs, batch size 16, actor
 447 learning rate 5×10^{-7} with cosine warmup and 3% warmup ratio, critic learning rate 9×10^{-6} , KL
 448 coefficient 0.0, rollout temperature 1.0, and 8 votes per prompt. Validation uses temperature 0.6 with
 449 8 samples. Models are trained for 1 epoch with checkpoint selection based on the best validation
 450 accuracy, maj@8.

451 **Data selection configuration.** For the neuron-guided data selection used by N-OPSD, we apply the
 452 same rule across all domains: *bot20*, which retains the bottom-20% of the source pool by zero-shot
 453 sensitivity, where $s(x) = |\mathcal{N}_{0\text{-shot}}(x)|$ is the number of unique neurons activated without context. The
 454 selected pool is then paired with $K = 10$ Jaccard-nearest few-shot demonstrations during training;
 455 see §4.3. No ground-truth labels are used at any point during selection.

456 B Methodology

457 B.1 Algorithm

Algorithm 1: Neuron-OPSD

Input: Unlabeled pool \mathcal{D} ; LLM π_θ ; percentile b ; shots K ; training steps T ; EMA rate ρ ;
 learning rate η .

Output: Updated model π_θ .

// Stage 1: neuron activation extraction (§4.1)

foreach $x \in \mathcal{D}$ **do**

 | Generate $y \sim \pi_\theta(\cdot | x)$ and score neurons via Eq. 3
 | Form activation set $\mathcal{N}(x)$ and neuron count $s(x) = |\mathcal{N}(x)|$

// Stage 2: high-consensus selection (§4.2)

$\tau_b \leftarrow b$ -th percentile of $\{s(x)\}_{x \in \mathcal{D}}$

$\mathcal{S} \leftarrow \{x \in \mathcal{D} \mid s(x) \leq \tau_b\}$

// Stage 3: context construction (§4.3)

Precompute $J(x_i, x_j)$ for all pairs in \mathcal{D} via Eq. 5

foreach $x \in \mathcal{S}$ **do**

 | $\mathcal{C}_K(x) \leftarrow K$ -nearest neighbors of x in $\mathcal{D} \setminus \{x\}$ under J

// Stage 4: self-training (§4.4)

Initialise EMA teacher $\theta' \leftarrow \theta$

for $t = 1, \dots, T$ **do**

 | Sample batch $B \subset \mathcal{S}$

 | Generate student rollouts $y \sim \pi_\theta(\cdot | x)$ for $x \in B$

 | Compute \mathcal{L}_{OPD} between $\pi_\theta(\cdot | x)$ and $\pi_{\theta'}(\cdot | x, \mathcal{C}_K(x))$ via Eq. 10

 | $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_{\text{OPD}}$

 | $\theta' \leftarrow (1 - \rho) \theta' + \rho \theta$

// EMA update

return π_θ

458 C Experiments

459 C.1 Evaluation

460 All evaluations sample 8 responses per question with temperature 0.6 using vLLM. We report
 461 majority-vote accuracy, maj@8, and Expected Calibration Error (ECE), where the confidence for
 462 each question is $c = \text{maxcount}/8$.

463 C.2 Case study: neuron-Jaccard retrieval on CHEM and PHYS.

464 To make the reasoning-diversity claim concrete, we run the same $K=10$ neuron-Jaccard retriever on
 465 two different queries within each of CHEM and PHYS., and inspect the retrieved few-shot contexts.

466 The contrast is not in the *overlap* between the two queries' demo sets, which is near-zero given the
 467 large sample pool, but in the *internal coherence* of each query's retrieved set.

468 **CHEM: retrieval locks onto query-specific reasoning modes.** Query A asks about iodine-131
 469 half-life calculation. Its retrieved set is dominated by radioactivity, isotope decay, molarity, mass
 470 spectrometry, and other per-sample formula-application chemistry tasks. Query B asks about catalytic
 471 oxidation depolymerisation. Its retrieved set contains ten demonstrations all concerning polymer
 472 breakdown: polymer recycling, controlled-radical-polymerisation reverse processes, dynamic-covalent
 473 gel degradation, silicone degradation, and lignin depolymerisation, among others. Each
 474 query's retrieved context is internally thematic and distinct from the other query's.

475 **PHYS.: retrieval returns generic grab-bags regardless of query.** Query A asks for the acceleration
 476 of a 15 kg object under a 10 N force. Its retrieved set is thematically unrelated: an ML-framework
 477 question, a stress-intensity-factor method, a heat-engine efficiency problem, a Material-Point-Method
 478 description, a specific-heat calculation, a quantum-heat-engine note, a diffuse-interface theory passage,
 479 and a PEDOT:PSS conductivity question. Query B asks about the Good Regulator Theorem. Its
 480 retrieved set is similarly scattered, spanning van der Waals force, ionizing radiation, cosmic-ray acceleration,
 481 resonance-frequency design, CMB temperature, magnetic materials, plasma electron density,
 482 and IR polarization. The two sets differ in specific questions but are structurally indistinguishable as
 483 generic physics expository or MCQ snippets.

484 **Interpretation.** In CHEM, samples traverse genuinely different reasoning paths, so neuron-level
 485 similarity can surface query-specific clusters. In PHYS., samples reuse largely the same reasoning
 486 neurons across distinct topics, so K -nearest neurons cannot differentiate queries and the retriever
 487 collapses toward generic contexts.

488 C.3 Relational Analysis of Neuron for Hallucination Detection

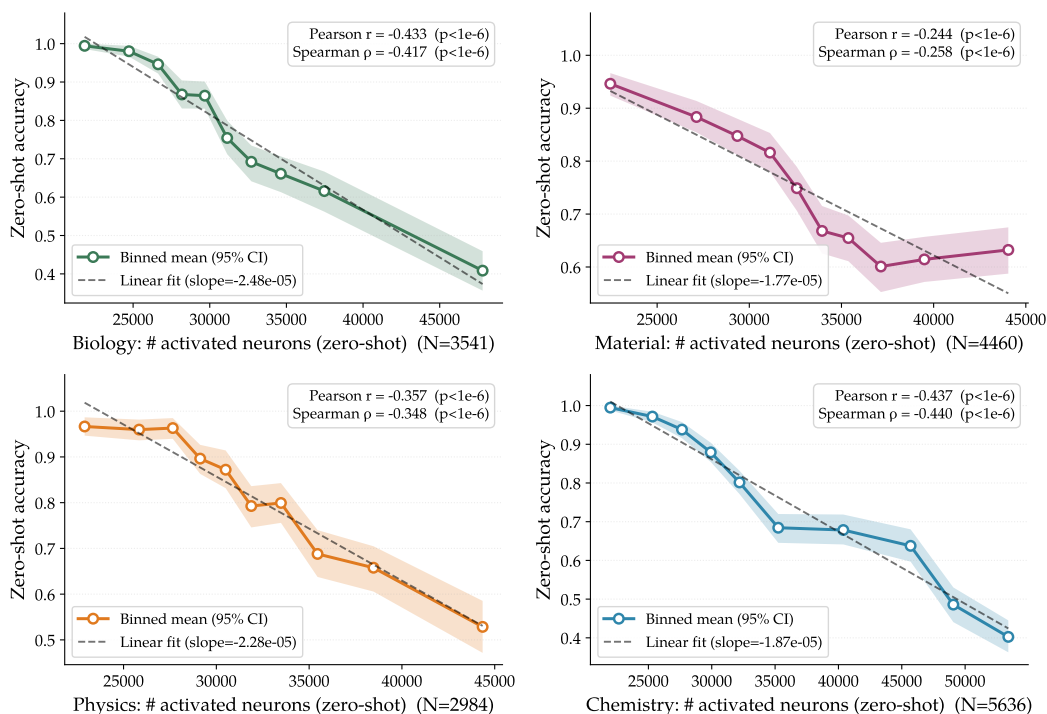


Figure 3: Relational regression results on SciKnowEval.

489 C.4 Additional Results

Model	SciKnowEval				Avg.	Edu.	MMLU-Pro	Avg.
	BIO	MAT.	PHYS.	CHEM				
Qwen3-4B	73.98	71.12	80.20	72.11	74.35	64.95	72.16	72.42
LMSI [Huang et al., 2023]								
LMSI _{BIO}	74.59 ^{+0.61}	71.22 ^{+0.10}	78.43 ^{-1.77}	68.60 ^{-3.51}	73.21 ^{-1.14}	66.00 ^{+1.05}	58.59 ^{-13.57}	69.57 ^{-2.85}
LMSI _{MAT.}	72.75 ^{-1.23}	69.08 ^{-2.04}	73.80 ^{-6.40}	69.32 ^{-2.79}	71.23 ^{-3.12}	60.32 ^{-4.63}	51.19 ^{-20.97}	66.08 ^{-6.34}
LMSI _{PHYS.}	73.70 ^{-0.28}	72.30 ^{+1.18}	78.54 ^{-1.66}	71.61 ^{-0.50}	74.04 ^{-0.31}	59.49 ^{-5.46}	58.01 ^{-14.15}	68.94 ^{-3.48}
LMSI _{CHEM}	74.56 ^{+0.58}	72.20 ^{+1.08}	81.40 ^{+1.20}	73.72 ^{+1.61}	75.47 ^{+1.12}	61.79 ^{-3.16}	65.99 ^{-6.17}	71.61 ^{-0.81}
LMSI _{EDU.}	62.85 ^{-11.13}	58.44 ^{-12.68}	64.98 ^{-15.22}	55.06 ^{-17.05}	60.34 ^{-14.02}	66.95 ^{+2.00}	49.37 ^{-22.79}	59.61 ^{-12.81}
LMSI _{MMLU}	74.27 ^{+0.29}	72.63 ^{+1.51}	81.06 ^{+0.86}	71.51 ^{-0.60}	74.87 ^{+0.52}	67.84 ^{+2.89}	68.60 ^{-3.56}	72.65 ^{+0.23}
LMSI _{Avg.}	72.12 ^{-1.86}	69.31 ^{-1.81}	76.37 ^{-3.83}	68.30 ^{-3.81}	71.53 ^{-2.83}	63.73 ^{-1.22}	58.62 ^{-13.53}	68.08 ^{-4.34}
TTRL [Zuo et al., 2025]								
TTRL _{BIO}	75.00 ^{+1.02}	73.03 ^{+1.92}	83.28 ^{+3.09}	72.85 ^{+0.74}	76.04 ^{+1.69}	65.31 ^{+0.36}	72.08 ^{-0.08}	73.59 ^{+1.17}
TTRL _{MAT.}	74.21 ^{+0.22}	73.26 ^{+2.14}	83.25 ^{+3.05}	72.03 ^{-0.08}	75.69 ^{+1.33}	64.92 ^{-0.03}	72.20 ^{+0.04}	73.31 ^{+0.89}
TTRL _{PHYS.}	74.14 ^{+0.16}	72.61 ^{+1.49}	82.98 ^{+2.79}	72.61 ^{+0.50}	75.58 ^{+1.23}	65.09 ^{+0.14}	71.94 ^{-0.22}	73.23 ^{+0.81}
TTRL _{CHEM}	74.94 ^{+0.95}	72.81 ^{+1.69}	83.40 ^{+3.20}	72.43 ^{+0.32}	75.89 ^{+1.54}	64.21 ^{-0.74}	71.92 ^{-0.24}	73.28 ^{+0.86}
TTRL _{EDU.}	74.56 ^{+0.57}	72.88 ^{+1.76}	83.55 ^{+3.35}	71.91 ^{-0.20}	75.72 ^{+1.37}	64.95 ^{+0.00}	69.43 ^{-2.73}	72.88 ^{+0.46}
TTRL _{MMLU}	74.94 ^{+0.95}	72.88 ^{+1.76}	83.89 ^{+3.69}	72.17 ^{+0.06}	75.97 ^{+1.62}	64.91 ^{-0.04}	72.04 ^{-0.12}	73.47 ^{+1.05}
TTRL _{Avg.}	74.63 ^{+0.65}	72.91 ^{+1.79}	83.39 ^{+3.19}	72.33 ^{+0.22}	75.82 ^{+1.46}	64.90 ^{-0.05}	71.60 ^{-0.56}	73.29 ^{+0.87}
Intuitior [Zhao et al., 2025b] ³								
Intuitior _{BIO}	76.21 ^{+2.22}	72.10 ^{+0.98}	82.23 ^{+2.03}	71.31 ^{-0.80}	75.46 ^{+1.11}	55.70 ^{-9.25}	66.17 ^{-5.99}	70.62 ^{-1.80}
Intuitior _{MAT.}	73.98 ^{+0.00}	73.41 ^{+2.29}	82.83 ^{+2.64}	71.67 ^{-0.44}	75.47 ^{+1.12}	63.82 ^{-1.13}	71.63 ^{-0.53}	72.89 ^{+0.47}
Intuitior _{PHYS.}	74.49 ^{+0.51}	72.53 ^{+1.41}	83.43 ^{+3.24}	72.07 ^{-0.04}	75.63 ^{+1.28}	54.81 ^{-10.14}	67.76 ^{-4.41}	70.85 ^{-1.57}
Intuitior _{CHEM}	74.37 ^{+0.38}	72.20 ^{+1.08}	83.21 ^{+3.01}	72.63 ^{+0.52}	75.60 ^{+1.25}	54.94 ^{-10.01}	69.10 ^{-3.06}	71.08 ^{-1.34}
Intuitior _{MMLU}	74.11 ^{+0.13}	73.39 ^{+2.27}	83.51 ^{+3.31}	72.57 ^{+0.46}	75.89 ^{+1.54}	62.92 ^{-2.03}	71.91 ^{-0.25}	73.07 ^{+0.65}
Intuitior _{Avg.}	74.63 ^{+0.65}	72.73 ^{+1.61}	83.04 ^{+2.85}	72.05 ^{-0.06}	75.61 ^{+1.26}	58.44 ^{-6.51}	69.32 ^{-2.85}	71.70 ^{-0.72}
Neuron-OPSD (ours)								
N-OPSD _{BIO}	74.02 ^{+0.03}	72.96 ^{+1.84}	82.87 ^{+2.67}	71.97 ^{-0.14}	75.45 ^{+1.10}	66.31 ^{+1.36}	69.46 ^{-2.71}	72.93 ^{+0.51}
N-OPSD _{MAT.}	73.83 ^{-0.16}	73.39 ^{+2.27}	83.13 ^{+2.94}	72.83 ^{+0.72}	75.79 ^{+1.44}	65.34 ^{+0.39}	71.22 ^{-0.94}	73.29 ^{+0.87}
N-OPSD _{PHYS.}	74.78 ^{+0.79}	73.71 ^{+2.60}	83.13 ^{+2.94}	71.93 ^{-0.18}	75.89 ^{+1.54}	65.20 ^{+0.25}	71.24 ^{-0.92}	73.33 ^{+0.91}
N-OPSD _{CHEM}	74.24 ^{+0.26}	73.69 ^{+2.57}	83.21 ^{+3.01}	72.03 ^{-0.08}	75.79 ^{+1.44}	65.45 ^{+0.50}	71.36 ^{-0.80}	73.33 ^{+0.91}
N-OPSD _{EDU.}	75.13 ^{+1.14}	73.14 ^{+2.02}	83.02 ^{+2.82}	72.85 ^{+0.74}	76.03 ^{+1.68}	72.19 ^{+7.24}	71.96 ^{-0.21}	74.71 ^{+2.29}
N-OPSD _{MMLU}	74.43 ^{+0.44}	73.46 ^{+2.34}	83.47 ^{+3.28}	73.03 ^{+0.92}	76.10 ^{+1.75}	64.33 ^{-0.62}	72.04 ^{-0.12}	73.46 ^{+1.04}
N-OPSD _{Avg.}	74.40 ^{+0.42}	73.39 ^{+2.27}	83.14 ^{+2.94}	72.44 ^{+0.33}	75.84 ^{+1.49}	66.47 ^{+1.52}	71.21 ^{-0.95}	73.51 ^{+1.09}

Table 7: Cross-domain evaluation on Qwen3-4B-Thinking-2507, performance reported in Avg@8. Each row reports a single model trained on one source domain and evaluated on all remaining test sets. Subscripts give the Δ vs. the untrained baseline.

Model	SciKnowEval					EDU.	MMLU-PRO	Avg.
	BIO	MAT.	PHYS.	CHEM	Avg.			
Qwen3-4B-think	0.195	0.213	0.095	0.173	0.169	0.246	0.184	0.184
TTRL [Zuo et al., 2025]								
TTRL _{BIO}	0.204 _{+0.009}	0.220 _{+0.007}	0.118 _{+0.023}	0.169 _{-0.004}	0.178 _{+0.009}	0.247 _{+0.000}	0.178 _{-0.006}	0.189 _{+0.005}
TTRL _{MAT.}	0.203 _{+0.005}	0.208 _{-0.005}	0.124 _{+0.029}	0.181 _{+0.005}	0.179 _{+0.010}	0.242 _{-0.005}	0.186 _{+0.002}	0.191 _{+0.006}
TTRL _{PHYS.}	0.214 _{+0.018}	0.215 _{+0.002}	0.123 _{+0.028}	0.166 _{-0.006}	0.180 _{+0.011}	0.239 _{-0.008}	0.188 _{+0.004}	0.191 _{+0.007}
TTRL _{CHEM}	0.204 _{+0.009}	0.223 _{+0.010}	0.118 _{+0.022}	0.192 _{+0.019}	0.184 _{+0.015}	0.251 _{+0.005}	0.197 _{+0.013}	0.198 _{+0.013}
TTRL _{EDU.}	0.199 _{+0.003}	0.215 _{+0.003}	0.111 _{+0.015}	0.164 _{-0.009}	0.172 _{+0.002}	0.248 _{+0.002}	0.202 _{+0.018}	0.190 _{+0.005}
TTRL _{MMLU}	0.200 _{+0.005}	0.224 _{+0.011}	0.114 _{+0.019}	0.169 _{-0.003}	0.177 _{+0.008}	0.243 _{-0.004}	0.198 _{+0.014}	0.191 _{+0.007}
TTRL _{Avg.}	0.204 _{+0.005}	0.218 _{+0.002}	0.118 _{+0.022}	0.174 _{+0.001}	0.178 _{+0.005}	0.245 _{-0.002}	0.191 _{+0.007}	0.192 _{+0.007}
Intuitior [Zhao et al., 2025b]								
Intuitior _{BIO}	0.208 _{+0.013}	0.236 _{+0.023}	0.141 _{+0.046}	0.210 _{+0.035}	0.199 _{+0.030}	0.312 _{+0.066}	0.229 _{+0.045}	0.223 _{+0.039}
Intuitior _{MAT.}	0.209 _{+0.014}	0.212 _{-0.001}	0.124 _{+0.029}	0.166 _{-0.007}	0.178 _{+0.009}	0.246 _{+0.000}	0.188 _{+0.004}	0.191 _{+0.007}
Intuitior _{PHYS.}	0.213 _{+0.018}	0.229 _{+0.012}	0.117 _{+0.021}	0.176 _{+0.004}	0.184 _{+0.015}	0.341 _{+0.095}	0.210 _{+0.026}	0.214 _{+0.030}
Intuitior _{CHEM}	0.213 _{+0.017}	0.225 _{+0.012}	0.130 _{+0.034}	0.198 _{+0.026}	0.191 _{+0.022}	0.319 _{+0.073}	0.234 _{+0.050}	0.220 _{+0.036}
Intuitior _{MMLU}	0.197 _{+0.002}	0.205 _{-0.008}	0.113 _{+0.018}	0.142 _{-0.031}	0.164 _{-0.005}	0.267 _{+0.021}	0.174 _{-0.010}	0.183 _{-0.001}
Intuitior _{Avg.}	0.208 _{+0.012}	0.221 _{+0.008}	0.125 _{+0.030}	0.179 _{+0.006}	0.183 _{+0.014}	0.297 _{+0.051}	0.207 _{+0.023}	0.206 _{+0.022}
Neuron-OPSD, ours								
N-OPSD _{BIO}	0.193 _{-0.003}	0.199 _{-0.014}	0.111 _{+0.015}	0.159 _{-0.013}	0.165 _{-0.004}	0.221 _{-0.025}	0.194 _{+0.010}	0.180 _{-0.005}
N-OPSD _{MAT.}	0.206 _{+0.011}	0.207 _{-0.006}	0.121 _{+0.026}	0.156 _{-0.017}	0.173 _{+0.004}	0.219 _{-0.028}	0.185 _{+0.001}	0.182 _{-0.002}
N-OPSD _{PHYS.}	0.196 _{+0.001}	0.198 _{-0.012}	0.114 _{+0.019}	0.170 _{-0.003}	0.170 _{+0.001}	0.217 _{-0.029}	0.189 _{+0.005}	0.181 _{-0.004}
N-OPSD _{CHEM}	0.198 _{+0.003}	0.215 _{+0.002}	0.117 _{+0.022}	0.171 _{-0.002}	0.175 _{+0.006}	0.225 _{-0.021}	0.180 _{-0.004}	0.184 _{+0.000}
N-OPSD _{EDU.}	0.195 _{-0.000}	0.218 _{+0.005}	0.122 _{+0.027}	0.171 _{-0.002}	0.177 _{+0.007}	0.191 _{-0.056}	0.179 _{-0.005}	0.179 _{-0.005}
N-OPSD _{MMLU}	0.206 _{+0.010}	0.215 _{+0.002}	0.113 _{+0.018}	0.158 _{-0.014}	0.173 _{+0.004}	0.239 _{-0.008}	0.180 _{-0.004}	0.185 _{+0.001}
N-OPSD _{Avg.}	0.199 _{+0.004}	0.209 _{-0.002}	0.116 _{+0.021}	0.164 _{-0.005}	0.172 _{+0.003}	0.219 _{-0.027}	0.185 _{+0.001}	0.182 _{-0.002}

Table 8: Expected Calibration Error (ECE), where lower is better, on the same test sets as Tab. 7. Each cell shows the absolute ECE with the delta vs the *Base* row in subscript, with green for lower ECE and better calibration and red for worse calibration. Gray cells mark the source domain. Shaded rows show the column-wise average across the trained models per method.