

DENSER \neq BETTER: LIMITS OF ON-POLICY SELF-DISTILLATION FOR CONTINUAL POST-TRAINING

Meng Wang¹, Haohan Zhao¹, Wenzhuo Liu¹, Lu Yang¹, Geng Liu¹, Haiyang Guo^{2,3},
Guo-Sen Xie⁴, Gaofeng Meng^{1,2,3}, Hongbin Liu^{1,2,3}, Fei Zhu^{1*}

¹Centre for Artificial Intelligence and Robotics, HKISI, CAS ²Institute of Automation, CAS
³University of Chinese Academy of Sciences ⁴Nanjing University of Science and Technology

ABSTRACT

Continual post-training enables foundation models to acquire new knowledge while preserving existing capabilities. Recent work suggests that on-policy learning can mitigate forgetting, with on-policy self-distillation emerging as a particularly attractive approach. In this work, we revisit this optimistic view through self-distillation policy optimization (SDPO). Our experiments show that SDPO can accelerate in-domain specialization when teacher signals are stable and well aligned, but it struggles to generalize to out-of-distribution scenarios. In continual post-training, SDPO exhibits stronger forgetting and can even collapse, whereas on-policy reinforcement learning methods such as GRPO adapt more conservatively and better preserve prior capabilities. Further analyses reveal that denser self-distillation induces larger drift in both parameter space and response space, and can amplify high-frequency formatting artifacts through a self-reinforcing teacher–student loop. These findings suggest that on-policy data alone is insufficient for continual learning. Dense self-distillation can accelerate specialization when teacher targets are stable and token-level supervision is reliable, but it should not be treated as a default stabilizer for continual post-training. Our code is available at <https://github.com/Moenupa/SDPO-CL>.

1 INTRODUCTION

Reinforcement learning with verifiable rewards has emerged as an effective post-training paradigm for large language models (LLMs) (Lambert et al., 2024; Guo et al., 2025a). However, practical systems often require continual post-training across domains and skills (Team, 2025; NVIDIA, 2025), raising a significant challenge: models must acquire new capabilities without forgetting prior ones. Existing work has suggested an optimistic picture for continual post-training: on-policy reinforcement fine-tuning (RFT) forgets less than supervised fine-tuning (SFT) (Lai et al., 2025; Chen et al., 2025; Shenfeld et al., 2025; Zhang et al., 2025b), and this benefit has often been attributed to the use of on-policy data (Chen et al., 2025; Shenfeld et al., 2025; Zhu et al., 2025).

Recently, on-policy distillation has emerged as an attractive alternative to sparse reward-based reinforcement learning. By converting model-generated rollouts into dense token-level supervision, these methods promise better credit assignment, higher sample efficiency, and reduced dependence on external teachers or explicit reward models (Agarwal et al., 2024; Gu et al., 2024; Thinking Machines Lab, 2025). A particularly appealing variant is on-policy self-distillation, where the model itself acts as a teacher by conditioning on demonstrations, feedback, solutions, or other auxiliary context, and the resulting behavior is distilled back into the same model. For example, self-distillation fine-tuning (SDFT) uses a demonstration-conditioned model as its own teacher and demonstrates that such on-policy distillation can acquire new skills while preserving prior capabilities (Shenfeld et al., 2026). Self-distillation policy optimization (SDPO) similarly uses a more-informed, feedback-conditioned teacher to provide dense token-level supervision to the student policy (Hübötter et al., 2026). Together, these works suggest that on-policy self-distillation may be a practical recipe for

*Corresponding author. Email: fei.zhu@cair-cas.org.hk

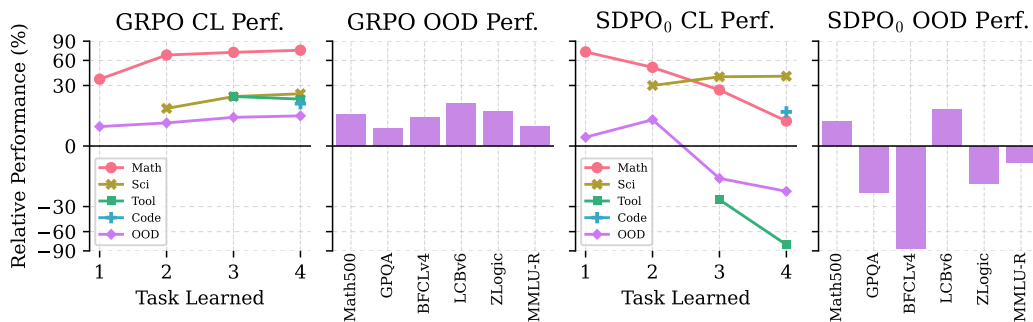


Figure 1: GRPO versus SDPO relative performance to base model (%) across sequentially learned tasks (1–4) and out-of-distribution (OOD) benchmarks. The line plots track performance retention across domains (MATH, SCIENCE, TOOLUSE, CODING, and macro OOD average) as new tasks are acquired. The bar charts break down final relative performance across specific OOD benchmarks.

continual learning, combining the on-policy nature of RFT that mitigates forgetting with the efficiency of dense supervised updates for learning new knowledge.

However, this view conflates two distinct ingredients: where the data come from and which objective is used to update the model. On-policy sampling determines the histories on which the model is trained, but the training objective determines which continuation distribution is reinforced at each history. This distinction is especially important for dense self-distillation. When every token becomes a training target, teacher drift, noisy rationales, formatting conventions, and domain-specific artifacts can be reinforced as repeatedly as useful behavior (Li et al., 2026; Kim et al., 2026). Thus, on-policy self-distillation may avoid the off-policy mismatch of SFT, but it does not automatically inherit the conservative update bias of sequence-level on-policy RFT.

This paper studies whether on-policy self-distillation enables reliable continual post-training. We focus on self-distillation policy optimization (SDPO, Hübötter et al., 2026) and compare it with sequence-level reward optimization (GRPO, Shao et al., 2024), from two perspectives:

- What are the key factors that make on-policy self-distillation beneficial? (§ 3)
- Does SDPO outperform GRPO in generalization or continual learning? (§ 4)

To answer these questions, we conduct post-training experiments in both single-domain and multi-domain continual settings. We vary the supervision density by comparing GRPO, standard SDPO, and SDPO variants with additional teacher strategies (§ 3.1) or supervision tweaks (§ 3.2), and evaluate on in-distribution and general benchmarks to assess specialization, retention, and transfer (§ 4.1 and 4.2). Finally, we connect behavior to diagnostics—parameter and response drift (§ 5.1), collapse modes (§ 5.2), and theoretical analysis (§ 5.3) to interpret our observations.

Our results show that on-policy self-distillation offers a powerful yet fragile training signal (Figure 1). SDPO can substantially improve performance on the current training domain, but it also increases the risk of drift, interference, and even collapse. Our observations suggest that supervision density introduces a trade-off: it strengthens the local learning signal, but also increases sensitivity, domain mismatch, and accumulated artifacts. Our main contributions are:

- We separate two factors often conflated in continual post-training: on-policy data and the objective used to turn on-policy rollouts into updates. This lets us ask whether on-policy self-distillation inherits the retention benefits attributed to on-policy RL.
- We show that SDPO is a strong specialization method but a fragile continual-learning signal. Its benefits depend on teacher stability and token reliability; CoT and fast teacher updates can amplify noise and artifacts.
- Across single-domain and staged continual post-training, SDPO exhibits weaker retention than GRPO, with forgetting concentrated on neighboring-yet-misaligned tasks and accumulating across domains.
- Parameter drift, response drift, collapse analysis, and an excess-KL view diagnose how and why SDPO yields such forgetting behavior.

2 PRELIMINARIES

2.1 POST-TRAINING OF LARGE LANGUAGE MODELS

Reinforcement Learning for Post-Training. Reinforcement learning is widely used to post-train large language models with preference signals or reward functions derived from human feedback, automated evaluators, or verifiable correctness criteria (Stiennon et al., 2020; Ouyang et al., 2022; Ziegler et al., 2019; Bai et al., 2022). Recent methods, such as direct preference optimization (Rafailov et al., 2023), recast preference optimization as a supervised objective, while on-policy algorithms like group relative policy optimization (Shao et al., 2024) optimize relative advantages over sampled response groups, often improving stability and scalability in large-scale settings.

On-Policy Distillation. On-policy distillation (OPD) methods construct training targets from samples generated by the current policy or its improved variants, thereby converting reinforcement signals into supervised learning updates (Agarwal et al., 2024; Gu et al., 2024). Recent self-distillation methods remove the need for a separate teacher by running the same model under different information conditions: the student observes only the original prompt, while the teacher is augmented with privileged context c . The privileged signal can come from demonstrations or examples (Snell et al., 2022; Shenfeld et al., 2026), correct solutions or textual feedback (Hübötter et al., 2026), or revision contexts, hints, and task-level feedback transformations (Zhao et al., 2026; Ye et al., 2026).

Chain-of-Thought Reasoning. Chain-of-thought (CoT) prompting and training encourage LLMs to generate intermediate reasoning before final answers, often improving performance on arithmetic, symbolic, and commonsense tasks (Wei et al., 2022; Kojima et al., 2022). Under final-answer supervision, this lengthens the credit-assignment horizon; with step-level or process rewards, it offers denser feedback and improves interpretability (Uesato et al., 2022; Lightman et al., 2023). CoT also has drawbacks: it increases inference cost, may propagate errors across steps, and the generated rationales need not faithfully reflect the model’s actual decision process (Turpin et al., 2023).

Continual Post-Training. Continual learning (Guo et al., 2025b) studies how models incorporate new knowledge without catastrophic forgetting of acquired ones. Common approaches use regularization, replay, or reference-model constraints to mitigate catastrophic forgetting across training stages (Kirkpatrick et al., 2017; Zhu et al., 2021; Liu et al., 2025). A growing line of work shows that RFT can substantially mitigate forgetting in continual post-training (Lai et al., 2025; Shenfeld et al., 2025; Chen et al., 2025). For example, Lai et al. (2025) find that this benefit is not mainly due to explicit KL penalties or chain-of-thought prompting, but rather to the implicit regularization induced by reward-guided policy updates. Complementarily, Shenfeld et al. (2025) argue that on-line RL is biased toward KL-minimal solutions among the policies that solve a new task, thereby reducing the distributional shift from the base model that is predictive of forgetting. Similarly, Chen et al. (2025) identify the use of on-policy data as a key factor behind RL’s robustness.

2.2 GROUP RELATIVE POLICY OPTIMIZATION

Group relative policy optimization (GRPO, Shao et al., 2024) assigns rewards at the sequence level. Given a prompt $x \sim \mathcal{D}$, GRPO samples a group of G responses $\{y_i\}_{i=1}^G$ from the current policy, i.e., $y_i \sim \pi_{\theta_{\text{old}}}(\cdot | x)$, evaluates each response with a reward model to obtain rewards $\{r_i\}_{i=1}^G$, and computes group-relative advantages by normalizing rewards within the group, e.g., $A_i = (r_i - \text{mean}(\{r_j\}_{j=1}^G)) / (\text{std}(\{r_j\}_{j=1}^G) + \epsilon)$. Training optimizes a PPO-style clipped objective:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \frac{1}{G} \sum_{i=1}^G \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \min(\rho_{i,t} A_i, \text{clip}(\rho_{i,t}, 1-\epsilon, 1+\epsilon) A_i) - \lambda D_{\text{KL}}(\pi_{\theta}(\cdot | x) \parallel \pi_{\text{ref}}(\cdot | x)),$$

where $\rho_{i,t} = \pi_{\theta}(y_{i,t} | x, y_{i,<t}) / \pi_{\theta_{\text{old}}}(y_{i,t} | x, y_{i,<t})$ and π_{ref} is a reference policy used for KL regularization. In essence, GRPO is a critic-free policy optimization method that stabilizes reinforcement learning by using relative rewards among multiple sampled responses to the same prompt.

2.3 ON-POLICY SELF-DISTILLATION

Self-distillation policy optimization (SDPO, Hübötter et al., 2026) instead provides dense token-level rewards using the same model as both teacher and student. For each student-generated response

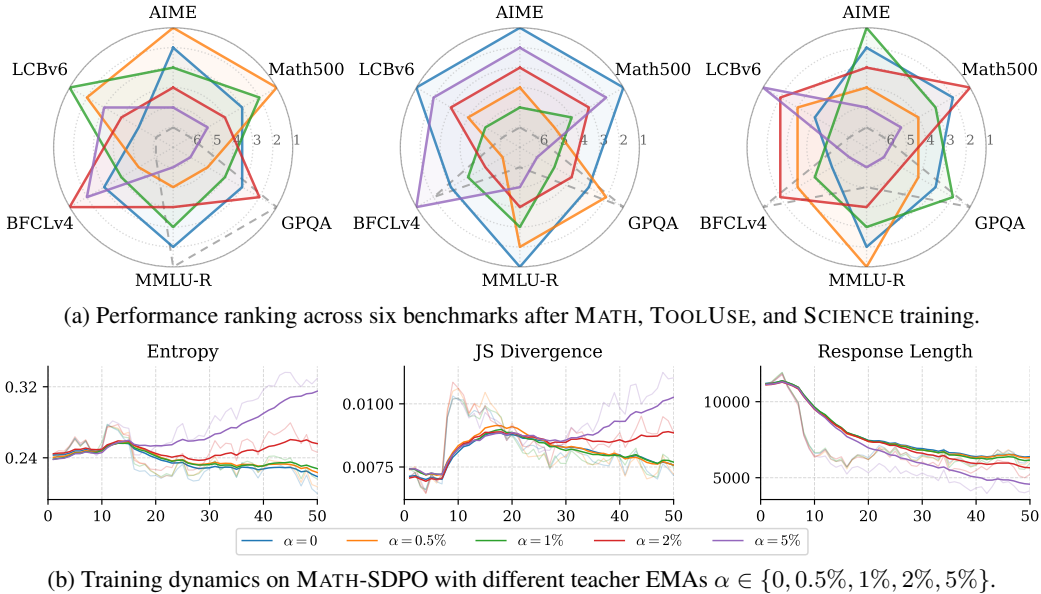


Figure 2: Effect of SDPO teacher EMA update rate α . No single α is globally optimal. **(a)** Ranking-based radar plots. A larger area is better. **(b)** Training dynamics during MATH-SDPO.

$y \sim \pi_\theta(\cdot | x)$, the student is trained to align with the more-informed teacher $\pi_\theta(\cdot | x, c)$ having richer context c (including solutions or other auxiliary signals). Training minimizes the generalized Jensen–Shannon divergence (JSD, Equation (A.2)) between student and teacher distributions:

$$\mathcal{L}_{\text{SDPO}}(\theta) = \sum_t D_{\text{JSD}}(\beta) (\pi_\theta(\cdot | x, y_{<t}) \parallel \text{stopgrad}(\pi_\theta(\cdot | x, c, y_{<t}))).$$

This yields dense supervision: instead of a single reward for the whole sequence, SDPO encourages student distribution at each generated token to align with more informed teacher distribution.

3 TOKEN-LEVEL DENSE SUPERVISION IS POWERFUL BUT SENSITIVE

First, we claim that SDPO’s benefits depend on the quality of the dense signal. SDPO distills token logits from the teacher model, improving sample efficiency because every generated token contributes supervision. However, it also means that teacher noise, temporal mismatch, and formatting artifacts are reinforced many times within a single response. We therefore begin by asking: *how do token-level reward stability (§ 3.1) and density (§ 3.2) affect self-distillation?*

3.1 TOKEN-LEVEL SUPERVISION REQUIRES FRESH BUT TEMPORALLY STABLE TEACHERS

A natural intuition in SDPO is that the teacher should be updated frequently: a fresher teacher should better reflect the student’s latest improvements and therefore provide stronger token-level supervision. To test this intuition, we sweep the EMA update rate $\alpha \in \{0, 0.5\%, 1\%, 2\%, 5\%\}$:

$$\Theta_{\text{teacher}}^{(k)} = (1 - \alpha)\Theta_{\text{teacher}}^{(k-1)} + \alpha\Theta_{\text{student}}^{(k)}.$$

Surprisingly, freshness alone does not predict performance. As shown in Figure 2a, no single α is uniformly optimal: in MATH-SDPO, $\alpha = 1\%$ gives the best overall profile, while $\alpha = 5\%$ performs worst; in TOOLUSE training, the frozen teacher is best, although $\alpha = 5\%$ is competitive. This raises a question: *if a larger EMA rate makes the teacher fresher, why can it hurt SDPO so severely?*

We argue that EMA entangles two different properties of the teacher: *freshness* and *temporal stability*. A larger α makes the teacher weights closer to the current student, but it also makes the teacher a more rapidly moving target. This is particularly risky for dense distillation: because SDPO applies supervision at every generated token, small teacher fluctuations can be reinforced repeatedly within

Table 1: MATH-SDPO performance (Acc@8). Restart-and-freeze offers a middle ground: it periodically re-initializes the teacher but prevents per-step noise from being distilled repeatedly.

Teacher Strategy	AIME	Math500	GPQA	ZLogic	MMLU-R	Avg
–	32.71	73.85	62.94	81.62	86.02	67.43
SDPO (EMA $\alpha = 0$)	56.42	84.29	58.72	74.50	85.47	71.88
StableSDPO ($\gamma = 1$)	57.50	85.05	60.86	74.68	86.60	72.94 ^{+1.06}
SDPO (EMA $\alpha = 5\%$)	34.38	78.97	55.87	73.68	83.82	65.34
StableSDPO ($\gamma = 40\%$)	55.00	84.62	59.03	73.50	87.05	71.84 ^{+6.50}

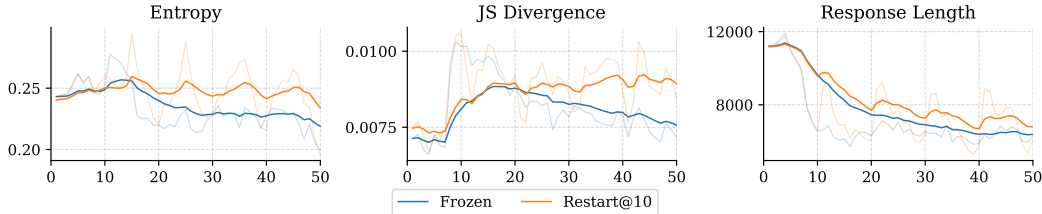


Figure 3: Training dynamics during MATH-SDPO with teacher *frozen* versus 10-step-hard-restart.

a single response. The training curves in Figure 2b support this view. Large EMA rates, especially $\alpha = 2\%$ and $\alpha = 5\%$, observe late-stage increases in entropy and JS divergence, indicating growing teacher–student mismatch and less reliable token targets. Smaller or frozen teachers avoid this instability, but they introduce the opposite problem: the teacher can become stale and provide a weakening optimization signal as the student evolves.

To separate freshness from per-step teacher dynamics, we introduce a simple *restart-and-freeze* strategy named StableSDPO. Instead of updating the teacher continuously, we periodically refresh it from the current student and then keep it frozen within the next interval:

$$\Theta_{\text{teacher}} \leftarrow (1 - \gamma)\Theta_{\text{teacher}}^{\text{old}} + \gamma\Theta_{\text{student}}^{(k)}, \quad 0 < \gamma \leq 1.$$

This design preserves teacher freshness at refresh points while removing the step-to-step volatility of EMA updates. For comparison with EMA, we match the effective k -step refresh amount by $\gamma \approx 1 - (1 - \alpha)^k$; for example, $\alpha = 5\%$ over $k = 10$ steps corresponds to $\gamma \approx 0.40$.

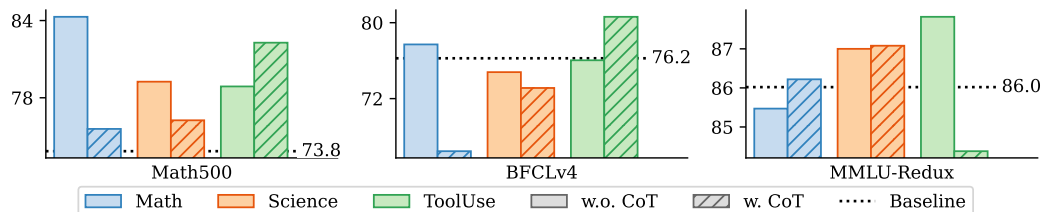
The results in Table 1 strongly support this decomposition. Replacing continuous EMA with matched restart-and-freeze substantially repairs the destructive behavior of $\alpha = 5\%$: on AIME, performance improves from 34.38% to 55.00%, and the average score improves from 65.34% to 71.84%, a gain of 6.50 points. Importantly, StableSDPO also improves over the frozen teacher: hard restart with $\gamma = 1$ outperforms the frozen teacher on all five benchmarks, improving the average score from 71.88% to 72.94%. Thus, the problem is not that freshness is harmful. Rather, freshness is useful only when it is delivered through a stable teacher signal. The training dynamics in Figure 3 further illustrate this point. The frozen teacher produces smooth curves, but its entropy and JS divergence gradually decline, suggesting that the teacher becomes increasingly outdated. In contrast, restart-and-freeze produces bounded oscillations aligned with refresh events: each restart renews the teacher signal, and avoids noisy per-step target changes. The JS divergence remains bounded, with neither late-stage blow-up nor collapse.

Takeaways 1

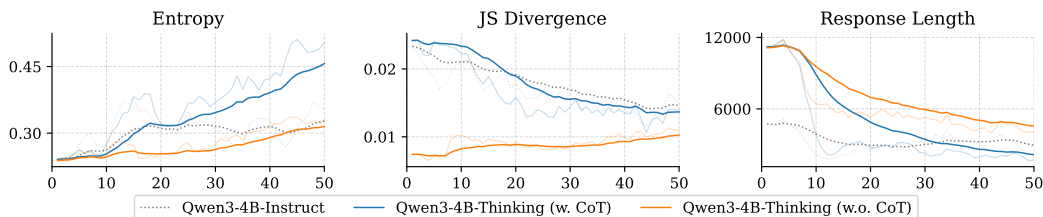
Teacher design is central to on-policy self-distillation, with a *stability–freshness* dilemma: a frozen teacher is stable but stale in the long term; a fast EMA teacher is fresh but can amplify online noise. Periodic refresh-and-freeze provides a more reliable dense target.

3.2 CoT DISTILLATION: MORE SUPERVISION IS NOT ALWAYS BETTER

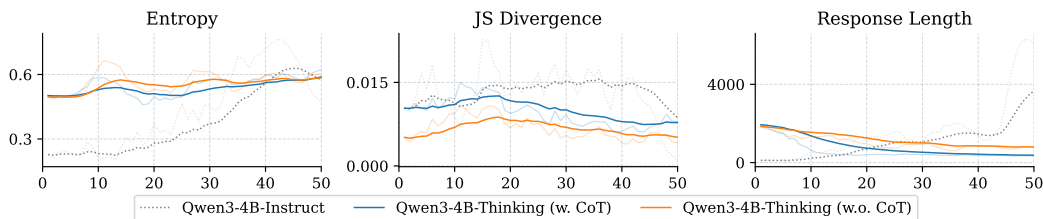
We next consider the supervision density of SDPO, by additionally distilling Chain-of-Thought tokens. This is a stronger test of the dense-supervision hypothesis: if more token-level targets were



(a) SDPO Performance (Acc@8) trained on MATH, SCIENCE, and TOOLUSE use against baseline (before training).



(b) SDPO training dynamics on MATH (which requires long CoT).



(c) SDPO training dynamics on TOOLUSE (which requires short CoT).

Figure 4: With and without CoT SDPO Performance (Acc@8) and training dynamics.

uniformly beneficial, CoT distillation should consistently improve performance. Figure 4a shows the opposite: CoT helps TOOLUSE training but hurts MATH and SCIENCE training.

We attribute this contrast to the reliability of intermediate tokens. CoT distillation does not merely add more supervision; it changes the target from answer-producing behavior to a particular reasoning trajectory. In long-form MATH and SCIENCE tasks, such trajectories are underdetermined and often contain redundant self-checks, uncertain interpretations, formatting habits, or incorrect intermediate claims. These tokens are weakly tied to final correctness, so distilling them can dilute credit assignment and encourage imitation of reasoning artifacts. This is reflected in MATH training by longer responses, higher entropy, and larger JS divergence (Figure 4b).

By contrast, TOOLUSE traces are short and schema-constrained. CoT therefore acts as useful process supervision in TOOLUSE training, yielding stable dynamics and improved BFCLv4 performance (Figure 4c). Additional EMA ablations further show that long-CoT supervision becomes especially harmful when the teacher is updated rapidly (Table D.2), consistent with a confirmation-bias loop in which noisy reasoning patterns are copied into the teacher and reinforced again.

Takeaways 2

More supervision is useful only when reliable. CoT distillation benefits short and structured traces, but can harm long, noisy, or artifact-prone rationales.

4 SDPO TRADES RETENTION FOR SPECIALIZATION

The previous section reveals that SDPO’s dense signal is conditional: it helps only when the teacher is stable and the distilled tokens are reliable. We now ask whether these local benefits translate into broader training advantages. *Does token supervision improves out-of-distribution (OOD) generalization (§ 4.1), and mitigates forgetting in continual post-training (§ 4.2)?*

Table 2: Performance (Acc@8) trained on a single domain. Rows indicate training domain and method; columns indicate evaluation benchmarks. SDPO₀ and SDPO_{5%} denote SDPO with teacher EMA 0 and 5%, respectively. SDPO₀^{CoT} denotes SDPO₀ with CoT distilled.

Method	MATH		SCIENCE		TOOLUSE		CODING		LOGIC	KNOW.	
	AIME	Math500	SCIENCE	GPQA	ToolUSE	BFCLv4	CODING	LCBv6	ZLogic	MMLU-R	
–	32.71	73.85	50.54	62.94	48.90	76.25	56.39	46.73	81.62	86.02	
MATH	GRPO	44.67	76.84	51.33	62.83↓	50.48	78.33	60.73	52.93	80.82↓	85.74↓
	SDPO _{5%}	34.38	78.97	45.00↓	55.87↓	53.31	78.12	52.37↓	50.86	73.68↓	83.82↓
	SDPO ₀	56.42	84.29	47.36↓	58.72↓	52.06	77.71	57.58	48.52	74.50↓	85.47↓
SCIENCE	GRPO	35.00	73.47↓	59.05	62.50↓	48.71↓	75.00↓	55.53↓	45.90↓	87.40	87.50
	SDPO _{5%}	53.34	78.42	72.26	57.39↓	51.29	72.29↓	60.66	56.11	90.68	84.10↓
	SDPO ₀	55.63	79.25	69.46	61.62↓	50.55	74.79↓	60.01	51.94	89.80	87.00
TOOLUSE	GRPO	38.55	76.60	49.11↓	62.25↓	57.90	79.58	57.88	47.69	77.35↓	84.95↓
	SDPO _{5%}	47.71	78.57	46.55↓	57.39↓	52.76	77.71	54.94↓	55.90	84.02	86.25
	SDPO ₀	47.71	78.88	50.65	60.73↓	55.51	76.04↓	61.00	58.48	86.48	87.82
CODING	GRPO	33.54	73.70↓	51.96	62.63↓	50.37	75.42↓	57.84	46.81	81.90	85.52↓
	SDPO _{5%}	49.17	77.98	47.92↓	58.33↓	49.63	73.96↓	62.10	58.47	80.53↓	87.25
	SDPO ₀	49.59	77.50	50.00↓	60.29↓	51.47	79.17	64.18	60.12	82.58	86.65

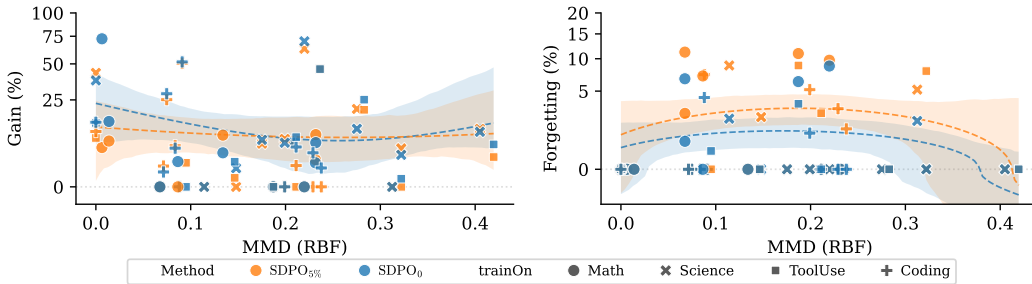


Figure 5: Relative Performance change as gain G and forgetting F with respect to Maximum Mean Discrepancy with a Radial Basis Function kernel.

4.1 RETHINKING OOD GENERALIZATION: INTERMEDIATE-DISTANCE INTERFERENCE

SDPO’s main advantage is specialization, not uniformly better generalization (Table 2). Under MATH training, SDPO₀ is the strongest math specialist, raising AIME from 32.71% to 56.42% and Math500 from 73.85% to 84.29%. Similarly, CODING-SDPO₀ improves LCBv6 by 13.39%. These results confirm that token-level supervision can strongly exploit source-like evaluation distributions.

However, this specialization does not translate into uniform retention. We decompose relative performance change into positive *gain* and negative *forgetting* (Equations (A.5) and (A.6)), and demonstrate that the SDPO-induced forgetting pattern is non-monotonic, as Figure 5 plots gain and forgetting against the distributional distance measured by Maximum Mean Discrepancy (MMD, Gretton et al., 2012). Source-like tasks (e.g., $MMD \lesssim 0.05$) and very distant tasks (e.g., $MMD \gtrsim 0.4$) mostly improve, while the largest drops concentrate in intermediate-distance tasks.

GPQA is the clearest intermediate-forget example: it is not the farthest benchmark, but consistently degrades after SDPO training from multiple source domains. We interpret this pattern as evidence that GPQA may lie in an interference region: it is close enough to be affected by source-domain updates, yet insufficiently aligned to benefit from them. More broadly, we observe SDPO degradation on MATH-to-SCIENCE, MATH-to-ZLogic, and SCIENCE-to-BFCLv4 (more cases are in Table D.4).

These results suggest that SDPO’s OOD failures are consistent with an intermediate-distance interference pattern, rather than a simple monotonic function of distributional distance. Distillation specializes the model toward the source-domain response manifold. When an OOD task is very close, this specialization transfers. When it is very far, the perturbation may not strongly affect it. The most vulnerable tasks are partially aligned OOD tasks: close enough to be changed, but

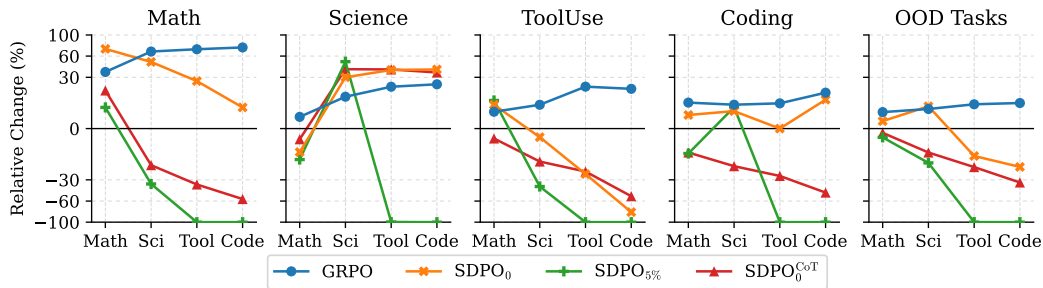


Figure 6: Normalized performance percentage change (y -axis) across continual training stages (x -axis). GRPO almost monotonically improves, while SDPO downshifts more often. Each subplot denotes one evaluation benchmark: MATH, SCIENCE, TOOLUSE, CODING, and out-of-domain tasks.

Table 3: Continual learning performance (Acc@8) of SDPO versus GRPO. SDPO often wins early (better current-stage), but GRPO wins over time (better last-stage).

Method	Stage	MATH	SCIENCE	TOOLUSE	CODING	LOGIC	KNOWLEDGE
Baseline	–	32.71	50.54	48.90	56.39	81.62	86.02
GRPO	Current	44.67 \uparrow 37%	56.39 \uparrow 12%	58.71 \uparrow 20%	–	–	–
	Last	57.29 \uparrow 75%	61.85 \uparrow 22%	57.72 \uparrow 18%	64.68 \uparrow 15%	89.52 \uparrow 10%	88.70 \uparrow 3%
SDPO ₀	Current	56.42 \uparrow 72%	65.70 \uparrow 30%	37.35 \downarrow 24%	–	–	–
	Last	34.38 \uparrow 5%	70.71 \uparrow 40%	9.93 \downarrow 80%	61.77 \uparrow 10%	72.15 \downarrow 12%	84.20 \downarrow 2%
SDPO ₀ ^{CoT}	Current	38.13 \uparrow 17%	70.89 \uparrow 40%	38.60 \downarrow 21%	–	–	–
	Last	14.17 \downarrow 57%	68.69 \uparrow 36%	23.35 \downarrow 52%	30.01 \downarrow 47%	59.95 \downarrow 27%	81.62 \downarrow 5%
SDPO _{5%}	Current	34.38 \uparrow 5%	76.37 \uparrow 51%	0.00 \downarrow 100%	–	–	–
	Last	0.00 \downarrow 100%	0.00 \downarrow 100%	0.00 \downarrow 100%	0.00 \downarrow 100%	0.02 \downarrow 100%	0.18 \downarrow 100%

misaligned enough not to benefit. This non-monotonic pattern is consistent with prior theory of continual learning with linear models (Lin et al., 2023; Goldfarb et al., 2024), and our results show that a similar interference regime emerges in LLM continual post-training,

Takeaways 3

For OOD generalization, dense supervision improves source-aligned tasks, leaves distant tasks relatively unaffected or even improved, but most strongly interferes with partially aligned tasks that are close enough to be shifted and misaligned enough not to benefit.

4.2 CONTINUAL LEARNING EXPOSES THE ACCUMULATED COST OF DENSE SUPERVISION

Single-domain interference becomes more severe when updates are composed across domains. We therefore train models sequentially on MATH \rightarrow SCIENCE \rightarrow TOOLUSE \rightarrow CODING and evaluate after each stage. We report both *Current*, i.e., performance on the current domain immediately after training it, and *Last*, i.e., performance after the full training sequence.

The continual setting reverses the single-domain picture. SDPO₀ is a strong early specialist: it reaches 56.42% on MATH and 65.70% on SCIENCE, exceeding GRPO’s 44.67% and 56.39%. However, these gains do not survive later domain shifts: after the full sequence, SDPO₀ drops to 34.38% on MATH and 9.93% on TOOLUSE, and also falls below the base model on LOGIC and KNOWLEDGE. SDPO₀^{CoT} degrades even earlier, consistent with the instability of long-CoT token targets. The most severe case is SDPO_{5%}, which collapses after the third stage and remains near zero. By contrast, GRPO adapts more conservatively but retains capabilities more reliably, achieving final improvements over the base model across all evaluated categories (Figure 6 and table 3).

These results argue that single-domain specialization is not evidence of continual learning. SDPO’s dense token-level updates can accelerate local adaptation, but the same updates also accumulate

drift, interference, and artifacts across domain shifts. Sequence-level GRPO is less aggressive, but its more selective reward signal better supports long-horizon retention.

Takeaways 4

For continual post-training, dense token-level supervision improves local adaptation but is not inherently retention-preserving; when composed across tasks, it can accumulate drift and interference that lead to forgetting.

5 UNDERSTANDING SDPO IN CONTINUAL POST-TRAINING

5.1 DENSE SUPERVISION CAUSES LARGER PARAMETER AND RESPONSE DRIFT

The benchmarks report that SDPO adapts more aggressively than GRPO. To understand why, we look inside the model and ask whether SDPO’s greater behavioral changes correspond to larger parameter shifts. For each LLM weight matrix (e.g., self-attention V) at training step k : $\mathbf{W}^{(k)} \in \mathbb{R}^{m \times n}$, we compute the compact singular value decomposition:

$$\mathbf{W}^{(k)} = \mathbf{U}^{(k)} \mathbf{\Sigma}^{(k)} \mathbf{V}^{(k)\top},$$

where $\mathbf{U}^{(k)} \in \mathbb{R}^{m \times d}$, $\mathbf{V}^{(k)} \in \mathbb{R}^{n \times d}$, and $d = \min(m, n)$. The SVD decomposes each matrix into input-output directions and their associated spectral strengths, and we use it to quantify how training changes the geometry, spectral magnitude, subspace orientation, and dominant parameter regions of the model. Specifically, we calculate the following metrics between $\mathbf{W}^{(0)}$ and $\mathbf{W}^{(k)}$:

- *Normalized spectral shift*: ℓ_2 -normalized drift of the singular-value profile in Equation (A.8). Larger values indicate substantial strength changes in dominant directions.
- *Principal-angle rotation*: maximum rotation within top- k left singular subspaces, as shown in Equation (A.11). A small value implies better preservation of dominant subspaces.
- *Principal-mask overlap*: intersection-over-union of principal masks before and after training per Equation (A.12). High values mean dominant entries remain important.
- *Principal-mask update sparsity*: intersection-over-union of the initial principal mask and the changed-parameter mask in Equation (A.14). A high value demonstrates that training directly modifies structurally important parameters.

Figures 7a and 7b highlight a consistent drift hierarchy. GRPO changes the model least: its spectral shift and subspace rotation remain small, and its principal-mask overlap stays near 99.9%. SDPO variants move farther, with $\text{SDPO}_0^{\text{CoT}}$ producing the largest subspace change. This pattern supports an accumulated-drift explanation of continual-learning failure: SDPO produces notably larger spectral shifts and more principal-parameter modifications than GRPO. The ordering $\text{SDPO}_0^{\text{CoT}} > \text{SDPO}_{5\%} > \text{SDPO}_0 > \text{GRPO}$ matches the aggressiveness of the supervision signal: CoT increases token coverage, while fast EMA increases target non-stationarity. This supports the view that denser or less stable token-level targets drive stronger internal updates.

Parameter drift is mirrored in the model’s outputs. Response-embedding similarities remain tightly concentrated under GRPO, whereas SDPO shifts the distribution downward and develops heavier lower tails (Figures 7c and 7d, embedding computed by Qwen3-Embedding-4B, Zhang et al., 2025a). This suggests that SDPO does not merely improve selected answers; it reshapes the response manifold more strongly than sequence-level GRPO.

Takeaways 5

Faster adaptation comes with larger drift. SDPO reshapes parameters and responses more than GRPO, which helps specialization but hurts continual stability.

5.2 COLLAPSE AS ARTIFACT AMPLIFICATION UNDER DENSE SELF-DISTILLATION

Finally, we diagnose one extreme failure mode as a case of artifact amplification in self-distillation. $\text{SDPO}_{5\%}$ collapses with a visible generation artifact: endlessly repeated “boxed” tokens. This

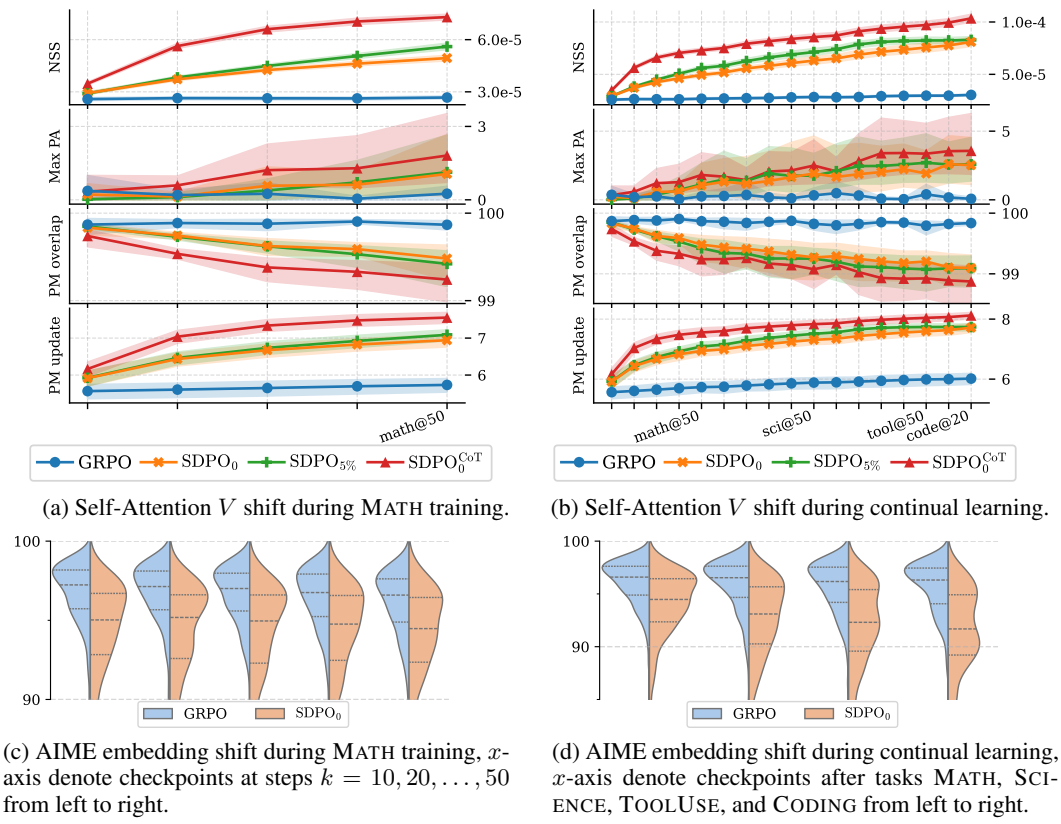


Figure 7: **(a, b)** Model parameter drift (y -axis) against training steps (x -axis). Metrics are averaged across all layers with 95% confidence interval. NSS, PA, PM denote normalized spectral shift, principal-angle, and principal-mask, respectively. **(c, d)** Cosine similarity (%) of rollout(@8) embeddings before and after GRPO (left half) versus SDPO₀ (right half) training.

behavior coincides with a change in final-answer formatting requirement at the third stage, where TOOLUSE differs from earlier MATH and SCIENCE stages (data samples in Appendix C.2). The collapsed model is over-amplifying the high-frequency pattern from earlier on-policy training stages, as JS divergence (Figure 8) plummets in the third stage, even to zero within the first few steps of TOOLUSE training. More experiments show that reordering data or masking tokens only partially mitigates this collapse (Tables D.5 and D.6).

This failure is consistent with the central diagnostic observed throughout our analyses. Because supervision is token-level, frequent tokens (e.g., formatting patterns) can receive repeated reinforcement. As the teacher is derived from the student, such artifacts can enter a feedback loop we call *confirmation bias*: the student emits the artifact, the teacher assigns it token-level support, and the next student update makes the artifact more likely. A frequently refreshed teacher or denser supervision can worsen this process and, in severe cases, lead to near-complete collapse.

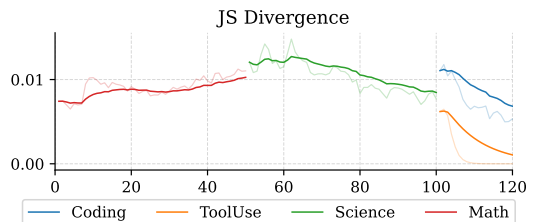


Figure 8: JS divergence during continual-SDPO_{5%} with TOOLUSE or CODING as third stage.

5.3 THEORETICAL VIEW ON FORGETTING: EXCESS KL DRIFT

Existing studies have revealed that on-policy RL is biased toward solutions that solve the new task while staying close to the base policy in KL (Shenfeld et al., 2025). However, SDPO does not

inherit this bias: its rollouts are on-policy, but its local target is a teacher distribution conditioned on augmented context, rather than the KL-minimal distribution induced by a sequence-level reward.

Let $x \sim \mathcal{D}$ be a prompt, y a response, and $S_x \subseteq \mathcal{Y}$ the set of successful responses under a binary sequence reward. Let $P_0(\cdot | x)$ be the sequence distribution induced by the base policy π_0 . For any success level $\alpha \in (0, 1)$, define the matched-success Razor policy

$$Q_{\alpha,x}^\dagger := \arg \min_{Q: Q(S_x)=\alpha} D_{\text{KL}}(Q \| P_0(\cdot | x)). \quad (1)$$

Let $q_{\alpha,x}^\dagger(\cdot | h_t)$ be its autoregressive token kernel at history $h_t = (x, y_{<t})$. This is the token distribution induced by the sequence-level reward: it preserves the base policy except for the reweighting needed to change the probability of eventual success. In the exact-success limit $\alpha \rightarrow 1$, $Q_{\alpha,x}^\dagger$ becomes the base-conditioned solution $P_0(\cdot | x, y \in S_x)$.

In contrast, SDPO distills a teacher that observes augmented context c , such as a solution, rationale, or auxiliary signal. Its local target is $T_c(\cdot | h_t) = \pi_T(\cdot | x, c, y_{<t})$. For a strict distillation divergence D , e.g., generalized JSD, the population token-level target is the teacher projection

$$B(\cdot | h_t) \in \arg \min_{u \in \Delta(\mathcal{V})} \mathbb{E}_{c|h_t} D(u, T_c(\cdot | h_t)). \quad (2)$$

When c is deterministic, this reduces to $B(\cdot | h_t) = T_c(\cdot | h_t)$. Let $P_B(\cdot | x)$ be the sequence distribution induced by this SDPO token target and set $\alpha_x = P_B(S_x | x)$.

Theorem 1 (Teacher projection induces excess KL drift). *Assume $0 < \alpha_x < 1$. Then*

$$\begin{aligned} & D_{\text{KL}}(P_B(\cdot | x) \| P_0(\cdot | x)) - D_{\text{KL}}(Q_{\alpha_x,x}^\dagger \| P_0(\cdot | x)) \\ &= \mathbb{E}_{y \sim P_B(\cdot | x)} \sum_{t=1}^{|y|} D_{\text{KL}}(B(\cdot | h_t) \| q_{\alpha_x,x}^\dagger(\cdot | h_t)) \geq 0. \end{aligned} \quad (3)$$

The inequality is strict whenever the SDPO teacher target differs from the sequence-reward Razor kernel on a set of histories with positive P_B -probability.

The no-excess-drift case is therefore restrictive: SDPO must satisfy $B(\cdot | h_t) = q_{\alpha_x,x}^\dagger(\cdot | h_t)$ at every visited history. This is not implied by teacher–student parameter sharing. Because the teacher conditions on augmented context c , it may select a particular reasoning trace, formatting pattern, or artifact even when many continuations are equally successful. Thus, SDPO’s on-policy sampling controls where supervision is applied, but the distribution selected at each token is the teacher projection, not the KL-minimal sequence-reward projection. Combined with the KL-forgetting law, the excess drift in Equation (3) predicts additional forgetting. Full proofs are in Appendix B.

Takeaways 6

On-policy data alone does not guarantee forgetting mitigation. SDPO trains on on-policy histories, but each token is projected toward a privileged teacher rather than a KL-minimal successful policy; this can amplify artifacts and induce excess drift across continual post-training.

6 CONCLUSION

This paper studied whether on-policy self-distillation improves continual post-training. Our results indicate that SDPO can be an effective specialization method when teacher signals are stable and aligned with the target domain. However, it is not uniformly beneficial for continual learning. SDPO’s token supervision amplifies both useful supervision and harmful artifacts, making it sensitive to teacher quality, update schedules, and token structure. Compared with GRPO, SDPO variants exhibit greater forgetting, instability, parameter drift, and response-manifold shifts.

The main lesson is not that on-policy self-distillation methods are ineffective, but they are not automatically safe. Token supervision can be a rapid specialization signal, but also a potentially dangerous signal: each token can reinforce useful signals and harmful artifacts. Sequence-level methods such as GRPO learn more slowly, but in our experiments, they preserve capabilities more reliably. Thus, the future of on-policy self-distillation for continual learning likely depends not on making the signal denser, but on safeguarding it through token weighting, teacher controls, or data balancing strategies to be more selective, stable, and artifact-aware.

REFERENCES

- Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos Garea, Matthieu Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-generated mistakes. In *International Conference on Learning Representations*, volume 2024, pp. 21246–21263, 2024.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Howard Chen, Noam Razin, Karthik Narasimhan, and Danqi Chen. Retaining by doing: The role of on-policy data in mitigating forgetting. *arXiv preprint arXiv:2510.18874*, 2025.
- Kehua Feng, Keyan Ding, Weijie Wang, Xiang Zhuang, Zeyuan Wang, Ming Qin, Yu Zhao, Jianhua Yao, Qiang Zhang, and Huajun Chen. Sciknoweval: Evaluating multi-level scientific knowledge of large language models, 2024.
- Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad Reza Ghasemi Madani, Claire Barale, Robert McHardy, Joshua Harris, Jean Kaddour, Emile van Krieken, and Pasquale Minervini. Are we done with mmlu?, 2024.
- Daniel Goldfarb, Itay Evron, Nir Weinberger, Daniel Soudry, and PAul HAnd. The joint effect of task similarity and overparameterization on catastrophic forgetting—an analytical model. In *The Twelfth International Conference on Learning Representations*, 2024.
- Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *The journal of machine learning research*, 13(1):723–773, 2012.
- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models. In *International Conference on Learning Representations*, volume 2024, pp. 32694–32717, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025a.
- Haiyang Guo, Fanhu Zeng, Fei Zhu, Jiayi Wang, Xukai Wang, Jingang Zhou, Hongbo Zhao, Wenzhuo Liu, Shijie Ma, Xu-Yao Zhang, et al. A comprehensive survey on continual learning in generative models. *arXiv e-prints*, pp. arXiv–2506, 2025b.
- Jonas Hübner, Frederike Lübeck, Lejs Behric, Anton Baumann, Marco Bagatella, Daniel Marta, Ido Hakimi, Idan Shenfeld, Thomas Kleine Büning, Carlos Guestrin, et al. Reinforcement learning via self-distillation. *arXiv preprint arXiv:2601.20802*, 2026.
- Naman Jain, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. In *International Conference on Learning Representations*, volume 2025, pp. 58791–58831, 2025.
- Jeonghye Kim, Xufang Luo, Minbeom Kim, Sangmook Lee, Dohyung Kim, Jiwon Jeon, Dongsheng Li, and Yuqing Yang. Why does self-distillation (sometimes) degrade the reasoning capability of llms? *arXiv preprint arXiv:2603.24472*, 2026.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.

-
- Song Lai, Haohan Zhao, Rong Feng, Changyi Ma, Wenzhuo Liu, Hongbo Zhao, Xi Lin, Dong Yi, Qingfu Zhang, Hongbin Liu, et al. Reinforcement fine-tuning naturally mitigates forgetting in continual post-training. *arXiv preprint arXiv:2507.05386*, 2025.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- Yaxuan Li, Yuxin Zuo, Bingxiang He, Jinqian Zhang, Chaojun Xiao, Cheng Qian, Tianyu Yu, Huanang Gao, Wenkai Yang, Zhiyuan Liu, and Ning Ding. Rethinking on-policy distillation of large language models: Phenomenology, mechanism, and recipe. *arXiv preprint arXiv:2604.13016*, 2026.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.
- Bill Yuchen Lin, Ronan Le Bras, Kyle Richardson, Ashish Sabharwal, Radha Poovendran, Peter Clark, and Yejin Choi. ZebraLogic: On the scaling limits of llms for logical reasoning. *arXiv preprint arXiv:2502.01100*, 2025.
- Sen Lin, Peizhong Ju, Yingbin Liang, and Ness Shroff. Theory on forgetting and generalization of continual learning. In *International Conference on Machine Learning*, pp. 21078–21100. PMLR, 2023.
- Wenzhuo Liu, Fei Zhu, Longhui Wei, and Qi Tian. C-clip: Multimodal continual learning for vision-language model. In *The Thirteenth International Conference on Learning Representations*, 2025.
- MAA. American invitational mathematics examination (aime), 2024–2025. URL <https://maa.org/maa-invitational-competitions/>.
- NVIDIA. Nvidia nemotron 3: Efficient and open intelligence, 2025. URL <https://arxiv.org/abs/2512.20856>. White Paper.
- Team Olmo, Allyson Ettinger, Amanda Bertsch, Bailey Kuehl, David Graham, David Heineman, Dirk Groeneveld, Faeze Brahman, Finbarr Timbers, Hamish Ivison, et al. Olmo 3. *arXiv preprint arXiv:2512.13961*, 2025.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- Shishir G. Patil, Huanzhi Mao, Charlie Cheng-Jie Ji, Fanjia Yan, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741, 2023.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Driani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. *arXiv preprint arXiv:2311.12022*, 2023.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Idan Shenfeld, Jyothish Pari, and Pulkit Agrawal. RL’s razor: Why online reinforcement learning forgets less. *arXiv preprint arXiv:2509.04259*, 2025.

-
- Idan Shenfeld, Mehul Damani, Jonas Hübotter, and Pulkit Agrawal. Self-distillation enables continual learning. *arXiv preprint arXiv:2601.19897*, 2026.
- Charlie Snell, Dan Klein, and Ruiqi Zhong. Learning by distilling context. *arXiv preprint arXiv:2209.15189*, 2022.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in neural information processing systems*, 33:3008–3021, 2020.
- Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, and Le Sun. Toolalpaca: Generalized tool learning for language models with 3000 simulated cases, 2023.
- Qwen Team. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Thinking Machines Lab. On-policy distillation. <https://thinkingmachines.ai/blog/on-policy-distillation/>, 2025. Blog post.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. *Advances in Neural Information Processing Systems*, 36:74952–74965, 2023.
- Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Tianzhu Ye, Li Dong, Xun Wu, Shaohan Huang, and Furu Wei. On-policy context distillation for language models. *arXiv preprint arXiv:2602.12275*, 2026.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *Advances in Neural Information Processing Systems*, 38:113222–113244, 2026.
- Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. Qwen3 embedding: Advancing text embedding and reranking through foundation models. *arXiv preprint arXiv:2506.05176*, 2025a.
- Zhihao Zhang, Qiaole Dong, Qi Zhang, Jun Zhao, Enyu Zhou, Zhiheng Xi, Senjie Jin, Xiaoran Fan, Yuhao Zhou, Mingqi Wu, et al. Why reinforcement fine-tuning enables mllms preserve prior knowledge better: A data perspective. *arXiv preprint arXiv:2506.23508*, 2025b.
- Siyan Zhao, Zhihui Xie, Mengchen Liu, Jing Huang, Guan Pang, Feiyu Chen, and Aditya Grover. Self-distilled reasoner: On-policy self-distillation for large language models. *arXiv preprint arXiv:2601.18734*, 2026.
- Fei Zhu, Xu-Yao Zhang, Chuang Wang, Fei Yin, and Cheng-Lin Liu. Prototype augmentation and self-supervision for incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5871–5880, 2021.
- Hanqing Zhu, Zhenyu Zhang, Hanxian Huang, DiJia Su, Zechun Liu, Jiawei Zhao, Igor Fedorov, Hamed Pirsiavash, Zhizhou Sha, Jinwon Lee, et al. The path not taken: RLvr provably learns off the principals. *arXiv preprint arXiv:2511.08567*, 2025.
- Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.

A MATH NOTATIONS

Kullback–Leibler (KL) divergence. For two discrete distributions $P(\mathcal{X})$ and $Q(\mathcal{X})$,

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right). \quad (\text{A.1})$$

Generalized Jensen–Shannon (JS) divergence. A weighted average of two KL divergences from P and Q to their mixture with $\beta \in (0, 1)$ ($\beta = 0.5$ recovers the standard symmetric JS divergence):

$$D_{\text{JSD}}(\beta)(P \parallel Q) = \beta D_{\text{KL}}(P \parallel \beta P + (1 - \beta)Q) + (1 - \beta) D_{\text{KL}}(Q \parallel \beta P + (1 - \beta)Q). \quad (\text{A.2})$$

Performance Change. We measure absolute (Δ) and relative (δ) change against the base model:

$$\Delta_{m,\text{train,eval}} = P_{m,\text{train,eval}} - P_{\text{base,eval}}, \quad (\text{A.3})$$

$$\delta_{m,\text{train,eval}} = \Delta_{m,\text{train,eval}} / P_{\text{base,eval}} \times 100\%, \quad (\text{A.4})$$

where m is the training method, train is the training domain, and eval is the evaluation benchmark. We define *gain* G and *forgetting* F as the positive and negative part of the relative change:

$$G_{m,\text{train,eval}} = \max(0, \delta_{m,\text{train,eval}}), \quad (\text{A.5})$$

$$F_{m,\text{train,eval}} = \max(0, -\delta_{m,\text{train,eval}}). \quad (\text{A.6})$$

Intersection-over-Union. For two masks A and B (e.g., predicted mask and ground-truth mask),

$$\text{IoU}(A, B) = |A \cap B| / |A \cup B|, \quad (\text{A.7})$$

where $|\cdot|$ denotes cardinality. The quantity lies in $[0, 1]$, with 1 for a perfect overlap.

Normalized spectral shift. It is defined as the ℓ_2 -normalized drift of the singular-value profile:

$$\text{NSS}(\mathbf{W}^{(0)}, \mathbf{W}^{(k)}) = \frac{\|\sigma(\mathbf{W}^{(0)}) - \sigma(\mathbf{W}^{(k)})\|_2}{\|\sigma(\mathbf{W}^{(0)})\|_2}. \quad (\text{A.8})$$

Principal-mask (PM) and Principal-angles (PAs). Given $\mathbf{W}_r^{(k)} = \mathbf{U}_r^{(k)} \Sigma_r^{(k)} \mathbf{V}_r^{(k)\top}$ as the rank- r approximation, PM is

$$\mathbf{M}_{\text{principal}}^{(k)}[i, j] = \mathbb{I} \left(\left| \mathbf{W}_r^{(k)}[i, j] \right| \geq \tau_\eta^{(k)} \right), \quad (\text{A.9})$$

where $\tau_\eta^{(k)}$ is a threshold controlled by the hyperparameter η . Principal angles are

$$\theta_i(\mathbf{W}^{(0)}, \mathbf{W}^{(k)}) = \arccos \sigma_i \left(\mathbf{U}_r^{(0)\top} \mathbf{U}_r^{(k)} \right), \quad i = 1, \dots, k, \quad (\text{A.10})$$

where $\sigma_i(\cdot)$ denotes the i -th singular value.

Maximum principal-angle rotation. Max PA of the top- r left singular vectors:

$$\text{maxPA}(\mathbf{W}^{(0)}, \mathbf{W}^{(k)}) = \max_{i \leq r} \theta_i(\mathbf{W}^{(0)}, \mathbf{W}^{(k)}). \quad (\text{A.11})$$

Principal-mask overlap. Intersection-over-Union between the initial and current PM masks:

$$\text{PMOverlap}(\mathbf{W}^{(0)}, \mathbf{W}^{(k)}) = \text{IoU} \left(\mathbf{M}_{\text{principal}}^{(0)}, \mathbf{M}_{\text{principal}}^{(k)} \right). \quad (\text{A.12})$$

Changed-parameter-mask. Per-parameter mask with change threshold $\eta > 0$:

$$\mathbf{M}_{\text{param}}^{(k)}[i, j] = \mathbb{I} \left(\left| \left(\mathbf{W}^{(k)} - \mathbf{W}^{(0)} \right)[i, j] \right| \geq \eta \right). \quad (\text{A.13})$$

Principal-mask update sparsity. Intersection-over-Union between the initial principal mask and the changed-parameter mask:

$$\text{PMUpdateSparsity}(\mathbf{W}^{(0)}, \mathbf{W}^{(k)}) = \text{IoU} \left(\mathbf{M}_{\text{principal}}^{(0)}, \mathbf{M}_{\text{param}}^{(k)} \right). \quad (\text{A.14})$$

B PROOFS

We provide the proofs for results in § 5.3. All distributions are over finite spaces. We assume that the base sequence distribution $P_0(\cdot | x)$ has full support. When the prompt x is fixed, we suppress it from notation.

Matched-success Razor policy. Let $S = S_x$ be the set of successful responses and let $z = P_0(S | x) \in (0, 1)$. For a target success level $\alpha \in (0, 1)$, define

$$\mathcal{Q}_\alpha = \{Q \in \Delta(\mathcal{Y}) : Q(S) = \alpha\}.$$

Lemma 2 (Closed form of the matched-success Razor policy). *The solution of*

$$Q_{\alpha,x}^\dagger = \arg \min_{Q \in \mathcal{Q}_\alpha} D_{\text{KL}}(Q \| P_0(\cdot | x))$$

is

$$Q_{\alpha,x}^\dagger(y) = P_0(y | x) \left[\frac{\alpha}{z} \mathbf{1}\{y \in S_x\} + \frac{1-\alpha}{1-z} \mathbf{1}\{y \notin S_x\} \right]. \quad (\text{B.1})$$

Equivalently,

$$Q_{\alpha,x}^\dagger = \alpha P_0(\cdot | x, S_x) + (1-\alpha) P_0(\cdot | x, S_x^c).$$

Proof. Fix x and write $p = P_0(\cdot | x)$. Let $p^+ = p(\cdot | S)$ and $p^- = p(\cdot | S^c)$. Any $Q \in \mathcal{Q}_\alpha$ can be decomposed as

$$Q = \alpha Q^+ + (1-\alpha) Q^-,$$

where $Q^+ = Q(\cdot | S)$ and $Q^- = Q(\cdot | S^c)$. Then

$$D_{\text{KL}}(Q \| p) = \alpha \log \frac{\alpha}{z} + (1-\alpha) \log \frac{1-\alpha}{1-z} + \alpha D_{\text{KL}}(Q^+ \| p^+) + (1-\alpha) D_{\text{KL}}(Q^- \| p^-).$$

The first two terms are constant over \mathcal{Q}_α . The remaining terms are nonnegative and vanish if and only if $Q^+ = p^+$ and $Q^- = p^-$. Thus the unique minimizer is

$$Q_{\alpha,x}^\dagger = \alpha p^+ + (1-\alpha) p^-,$$

which is exactly Equation (B.1). \square

Lemma 3 (Token kernel induced by the Razor policy). *Let $q_{\alpha,x}^\dagger(\cdot | h_t)$ be the autoregressive token kernel induced by $Q_{\alpha,x}^\dagger$ at history $h_t = (x, y_{<t})$. Define*

$$C_\alpha(h) = \mathbb{E}_{y \sim P_0(\cdot | h)} \left[\frac{\alpha}{z} \mathbf{1}\{y \in S_x\} + \frac{1-\alpha}{1-z} \mathbf{1}\{y \notin S_x\} \right]. \quad (\text{B.2})$$

Then

$$q_{\alpha,x}^\dagger(a | h_t) = \pi_0(a | h_t) \frac{C_\alpha(h_t a)}{C_\alpha(h_t)}. \quad (\text{B.3})$$

In the exact-success limit $\alpha \rightarrow 1$,

$$q_{1,x}^\dagger(a | h_t) = \pi_0(a | h_t) \frac{\Pr_{P_0}(y \in S_x | h_t a)}{\Pr_{P_0}(y \in S_x | h_t)}. \quad (\text{B.4})$$

Proof. By Theorem 2,

$$Q_{\alpha,x}^\dagger(y) = P_0(y | x) w_\alpha(y),$$

where

$$w_\alpha(y) = \frac{\alpha}{z} \mathbf{1}\{y \in S_x\} + \frac{1-\alpha}{1-z} \mathbf{1}\{y \notin S_x\}.$$

For a history $h_t = (x, y_{<t})$,

$$q_{\alpha,x}^\dagger(a | h_t) = \Pr_{Q_{\alpha,x}^\dagger}(y_t = a | h_t) = \frac{\sum_{y_{>t}} P_0(a, y_{>t} | h_t) w_\alpha(y)}{\sum_{\tilde{y}_{\geq t}} P_0(\tilde{y}_{\geq t} | h_t) w_\alpha(\tilde{y})}.$$

The numerator equals

$$\pi_0(a | h_t) \mathbb{E}_{y \sim P_0(\cdot | h_t a)}[w_\alpha(y)] = \pi_0(a | h_t) C_\alpha(h_t a),$$

and the denominator equals $C_\alpha(h_t)$. This proves Equation (B.3). Letting $\alpha \rightarrow 1$ gives Equation (B.4). \square

Teacher projection and excess reverse KL. For a strict distillation divergence D , the population SDPO token target is

$$B(\cdot | h_t) \in \arg \min_{u \in \Delta(\mathcal{V})} \mathbb{E}_{c|h_t} D(u, T_c(\cdot | h_t)).$$

When c is deterministic and D is generalized JSD, the minimizer is $B(\cdot | h_t) = T_c(\cdot | h_t)$ because generalized JSD is nonnegative and equals zero if and only if its two arguments are identical.

Let $P_B(\cdot | x)$ be the sequence distribution induced by $B(\cdot | h_t)$, and let

$$\alpha = P_B(S_x | x).$$

The corresponding matched-success Razor policy is $Q_{\alpha,x}^\dagger$.

Proof of Theorem 1. Fix x and suppress x from notation. Let $p = P_0(\cdot | x)$, $S = S_x$, $z = p(S)$, and $Q^\dagger = Q_{\alpha,x}^\dagger$. By Theorem 2,

$$Q^\dagger(y) = p(y)w_\alpha(y),$$

where $w_\alpha(y) = \alpha/z$ on S and $w_\alpha(y) = (1 - \alpha)/(1 - z)$ on S^c .

Since $P_B(S) = Q^\dagger(S) = \alpha$,

$$\sum_y P_B(y) \log w_\alpha(y) = \sum_y Q^\dagger(y) \log w_\alpha(y). \quad (\text{B.5})$$

Therefore,

$$\begin{aligned} D_{\text{KL}}(P_B \| p) - D_{\text{KL}}(Q^\dagger \| p) &= \sum_y P_B(y) \log \frac{P_B(y)}{p(y)} - \sum_y Q^\dagger(y) \log \frac{Q^\dagger(y)}{p(y)} \\ &= \sum_y P_B(y) \log \frac{P_B(y)}{p(y)} - \sum_y Q^\dagger(y) \log w_\alpha(y) \\ &= \sum_y P_B(y) \log \frac{P_B(y)}{p(y)w_\alpha(y)} \\ &= D_{\text{KL}}(P_B \| Q^\dagger), \end{aligned} \quad (\text{B.6})$$

where we used Equation (B.5) in the third line.

Now write both sequence distributions autoregressively:

$$P_B(y | x) = \prod_{t=1}^H B(y_t | h_t), \quad Q^\dagger(y | x) = \prod_{t=1}^H q_{\alpha,x}^\dagger(y_t | h_t).$$

Then the autoregressive KL chain rule gives

$$\begin{aligned} D_{\text{KL}}(P_B \| Q^\dagger) &= \mathbb{E}_{y \sim P_B} \log \frac{P_B(y | x)}{Q^\dagger(y | x)} \\ &= \mathbb{E}_{y \sim P_B} \sum_{t=1}^H \log \frac{B(y_t | h_t)}{q_{\alpha,x}^\dagger(y_t | h_t)} \\ &= \mathbb{E}_{y \sim P_B} \sum_{t=1}^H D_{\text{KL}}(B(\cdot | h_t) \| q_{\alpha,x}^\dagger(\cdot | h_t)). \end{aligned} \quad (\text{B.7})$$

Combining Equation (B.6) and Equation (B.7) proves the theorem. The right-hand side is nonnegative, and it is strictly positive whenever the two token kernels differ on a set of histories with positive P_B -probability. \square

Forward KL counterpart. The previous theorem uses reverse KL because RL's Razor is naturally stated as a reverse-KL projection. The empirical forgetting law is often measured by forward KL from the base policy to the fine-tuned policy, so we also give the matched-success forward-KL analogue.

Proposition 4 (Forward-KL excess at matched success). *Fix x and let $p = P_0(\cdot | x)$. Let $S = S_x$, $z = p(S)$, and let $P_B(S) = \alpha \in (0, 1)$. Define*

$$p^+ = p(\cdot | S), \quad p^- = p(\cdot | S^c),$$

and

$$B^+ = P_B(\cdot | S), \quad B^- = P_B(\cdot | S^c).$$

Then

$$D_{\text{KL}}(p \| P_B) - D_{\text{KL}}(p \| Q_{\alpha,x}^\dagger) = zD_{\text{KL}}(p^+ \| B^+) + (1-z)D_{\text{KL}}(p^- \| B^-) \geq 0. \quad (\text{B.8})$$

Proof. Decompose P_B as

$$P_B = \alpha B^+ + (1-\alpha)B^-.$$

Then

$$D_{\text{KL}}(p \| P_B) = z \log \frac{z}{\alpha} + (1-z) \log \frac{1-z}{1-\alpha} + zD_{\text{KL}}(p^+ \| B^+) + (1-z)D_{\text{KL}}(p^- \| B^-).$$

For $Q_{\alpha,x}^\dagger = \alpha p^+ + (1-\alpha)p^-$, the two conditional KL terms vanish:

$$D_{\text{KL}}(p \| Q_{\alpha,x}^\dagger) = z \log \frac{z}{\alpha} + (1-z) \log \frac{1-z}{1-\alpha}.$$

Subtracting gives Equation (B.8). \square

Connection to forgetting. We use the KL-forgetting law as an empirical bridge between distributional drift and prior-task forgetting.

Assumption 5 (KL-forgetting law). *For policies matched in new-task success, let $F(\pi)$ denote the drop in prior-task performance. Suppose there exists a strictly increasing function ω with $\omega(0) = 0$ such that*

$$\mathcal{K}_{\text{fwd}}(\pi_1) - \mathcal{K}_{\text{fwd}}(\pi_2) \geq \varepsilon \implies F(\pi_1) - F(\pi_2) \geq \omega(\varepsilon),$$

where

$$\mathcal{K}_{\text{fwd}}(\pi) = \mathbb{E}_{x \sim \mathcal{D}} D_{\text{KL}}(P_0(\cdot | x) \| P_\pi(\cdot | x)).$$

Corollary 6 (Teacher mismatch implies additional forgetting). *Let P_B be the SDPO teacher-projected sequence distribution and let $Q_{\alpha,x}^\dagger$ be the matched-success Razor policy with $\alpha = P_B(S_x | x)$. Define*

$$\begin{aligned} \Delta_{\text{fwd}} = \mathbb{E}_{x \sim \mathcal{D}} & \left[z_x D_{\text{KL}}(P_0(\cdot | x, S_x) \| P_B(\cdot | x, S_x)) \right. \\ & \left. + (1-z_x) D_{\text{KL}}(P_0(\cdot | x, S_x^c) \| P_B(\cdot | x, S_x^c)) \right]. \end{aligned}$$

If $\Delta_{\text{fwd}} > 0$, then under Theorem 5,

$$F(P_B) - F(Q_\alpha^\dagger) \geq \omega(\Delta_{\text{fwd}}) > 0.$$

Proof. Taking expectation over $x \sim \mathcal{D}$ in Theorem 4 gives

$$\mathcal{K}_{\text{fwd}}(P_B) - \mathcal{K}_{\text{fwd}}(Q_\alpha^\dagger) = \Delta_{\text{fwd}}.$$

The result follows directly from Theorem 5. \square

C TRAINING DETAILS

C.1 MODELS AND DATASETS

We adopt Qwen3 series (Qwen3-4B-Instruct-2507 and Qwen3-4B-Thinking-2507, Team, 2025), Olmo-3 series (Olmo-3-7B-Instruct-DPO and Olmo-3-7B-Think-DPO, Olmo et al., 2025) as base models¹ and train² on 4 diverse tasks—MATH, SCIENCE, TOOLUSE, and CODING, with in-distribution and out-of-distribution³ test sets to evaluate model performance measured in averaged accuracy (i.e., Acc@k). Data sources and training hyperparameters (learning rates according to Hübötter et al., 2026) are disclosed in Tables C.1 and C.2, respectively. We visualize the embeddings of each benchmark using Qwen3-Embedding-4B (Zhang et al., 2025a) in Figure C.1.

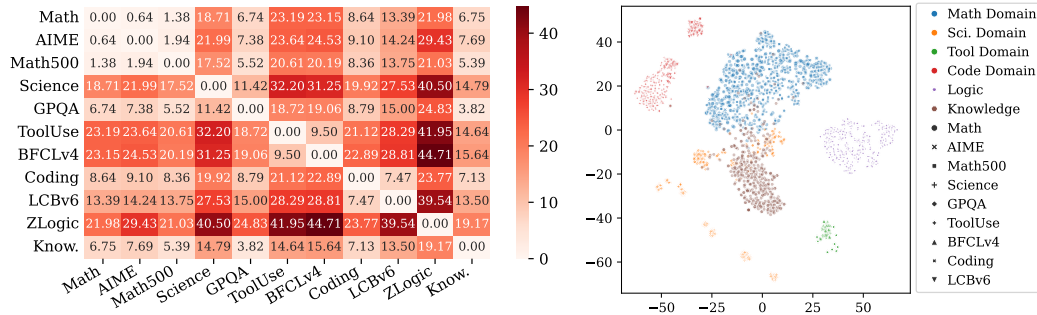
¹For Instruct models, SDPO and SDPO^{Cot} are equivalent due to absence of think tags.

²Training framework from SDPO official implementation <https://github.com/lasgroup/SDPO>.

³In-distribution and out-of-distribution are pragmatic labels, not strict measures of inter-dataset distance.

Table C.1: Data sources for training or evaluation purposes across domains. *AIME* denotes AIME24 and AIME25 combined. Abbreviations *ZLogic*, *MMLU-R* denote ZebraLogicBench, MMLU-Redux 2.0, respectively.

Domain	Training	Evaluation	
		In-Distribution	Out-Of-Distribution
MATH	DAPO-Math-17k (Yu et al., 2026)	AIME (MAA, 2024–2025)	Math500 (Lightman et al., 2023)
SCIENCE	SciKnowEval (Feng et al., 2024)		GPQA (Rein et al., 2023)
TOOLUSE	ToolAlpaca (Tang et al., 2023)		BFCLv4 (Patil et al., 2025)
CODING	Dolci-Think-RL-7B (Olmo et al., 2025)		LCBv6 (Jain et al., 2025)
LOGIC		–	ZLogic (Lin et al., 2025)
KNOWLEDGE		–	MMLU-R (Gema et al., 2024)



(a) Pairwise MMD (multiplier 10^{-2}) between benchmarks.

(b) Visualization using t-SNE.

Figure C.1: Embedding visualization of benchmarks by domain.

C.2 DATA SAMPLES

AIME

Every morning Aya goes for a s -kilometer-long walk and stops at a coffee shop afterwards. When she walks \rightarrow at a constant speed of s kilometers per hour, the walk takes her 4 hours, including t minutes spent in the coffee shop. When she walks $s+2$ kilometers per hour, the walk takes her 2 hours and 24 minutes, including t minutes spent in the coffee shop. Suppose Aya walks at $\frac{1}{2}$ kilometers per hour. Find the number of minutes the walk takes her, including the t minutes spent in \rightarrow the coffee shop.
Please reason step by step, and put your final answer within `\boxed{}`.

BFCLv4

Your task is to answer the user's question using available tools.
You have access to the following tools:
Name: Tools
Description: A set of functions/tools that you can use to answer the user's question.
Documentation:
math.triangle_area_heron: Calculates the area of a triangle using Heron's formula, given the lengths of \rightarrow its three sides.
Parameters: {"side1": "integer. Length of the first side of the triangle.", "side2": "integer. Length of \rightarrow the second side of the triangle.", "side3": "integer. Length of the third side of the triangle."}
Output: Successful response.
- Format: application/json
math.circle_area: Calculates the area of a circle given its radius.
Parameters: {"radius": "float. The radius of the circle."}
Output: Successful response.
- Format: application/json

Table C.2: Hyperparameters for SDPO and GRPO training with Qwen3-4B-Thinking-2507.

	Parameters	SDPO	GRPO
Data	Max. prompt length	2048	2048
	Max. response length	14336	14336
Batching	Question batch size	32	8
	Mini batch size	32	8
	Number of rollouts	8	8
Rollout	Inference engine	vLLM	vLLM
	Temperature	1.0	1.0
	Top- p	-1	-1
	Top- k	-1	-1
SDPO	Distillation Top- K	100	-
	Distillation divergence	Jensen-Shannon	-
	Teacher-EMA update rate (%)	$\alpha \in \{0, 0.5, 1, 2, 5\}$	-
	Rollout importance-sampling clip	2	-
PPO	KL coefficient (λ)	-	0.0
	ϵ -low	0.20	0.20
	ϵ -high	0.28	0.28
Training	Optimizer	AdamW	AdamW
	Learning rate	10^{-5}	10^{-6}
	Training steps	50 (20 for CODING)	50 (20 for CODING)
	Warmup steps	10	10
	Weight decay	0.01	0.01
	Gradient Clip Norm	1.0	1.0

```

math.triangle_area_base_height: Calculates the area of a triangle using the formula (1/2)base*height.
Parameters: {"base": "float. The base length of the triangle.", "height": "float. The height of the
↳ triangle."}
Output: Successful response.
- Format: application/json

Use the following format:
Thought: you should always think about what to do
Action: the action to take, should be one of the tool names.
Action Input: the input to the action, must be in JSON format. All of the action input must be realistic
↳ and from the user.

Begin!
Question: Calculate the area of a triangle, given the lengths of its three sides: 3, 4, and 5.

```

CODING

```

You are tasked with creating a function that retrieves the details of a specific assignment based on its
↳ ID. The function should take a string representing the assignment ID as input and return a dictionary
↳ containing the assignment details. For the purpose of this question, you can assume that the
↳ assignment details will be represented as a fixed dictionary for demonstration. Implement a function
↳ named 'get_assignment_details' that takes an assignment ID and returns the corresponding assignment
↳ details as a dictionary. If the assignment ID does not exist, return an empty dictionary. The
↳ assignment details for valid IDs are predefined as follows: '101': {'title': 'Math Homework',
↳ 'due_date': '2023-11-01'}, '102': {'title': 'Science Project', 'due_date': '2023-12-15'}, '103':
↳ {'title': 'History Essay', 'due_date': '2023-10-30'}.

```

GPQA

```

Two quantum states with energies E1 and E2 have a lifetime of 10^-9 sec and 10^-8 sec, respectively. We
↳ want to clearly distinguish these two energy levels. Which one of the following options could be their
↳ energy difference so that they can be clearly resolved?

```

- A. 10^{-8} eV
- B. 10^{-4} eV
- C. 10^{-11} eV
- D. 10^{-9} eV

Please reason step by step, and put your final answer within `\boxed{}`.

MATH

In triangle ABC , $\sin \angle A = \frac{4}{5}$ and $\angle A < 90^\circ$. Let D be a point outside $\triangle ABC$ such that $\angle BAD = \angle DAC$ and $\angle BDC = 90^\circ$. Suppose that $AD = 15$ and that $\frac{BD}{CD} = \frac{3}{2}$. If $AB + AC$ can be expressed in the form $\frac{a\sqrt{b}}{c}$ where a, b, c are pairwise relatively prime integers, find $a + b + c$. Please reason step by step, and put your final answer within `\boxed{}`.

Math500

Convert the point $(0, 3)$ in rectangular coordinates to polar coordinates. Enter your answer in the form (r, θ) , where $r > 0$ and $0 \leq \theta < 2\pi$. Please reason step by step, and put your final answer within `\boxed{}`.

MMLU-R

Determine whether the polynomial in $Z[x]$ satisfies an Eisenstein criterion for irreducibility over \mathbb{Q} . $x^2 - 12$
A. Yes, with $p=2$.
B. Yes, with $p=3$.
C. Yes, with $p=5$.
D. No.
Please reason step by step, and put your final answer within `\boxed{}`.

SCIENCE

Which of the following molecules is the only correct reactant used in the synthesis of the product "Nc1ccn(Cc2ccc(Cl)cc2C(F)(F)F)n1" ?

A: O=C1c2ccccc2C(=O)N1c1ccn(Cc2ccc(Cl)cc2C(F)(F)F)n1
B: O=C1c2ccccc2C(Cl)N1c1ccc2ccc(Cl)nc2n1
C: O=C1c2ccccc2C(=O)N1Cc1ccc(Oc2ccc(Cl)c(C(F)(F)F)c2)cc1
D: O=C1c2ccccc2C(=O)N1CC=C(Cc1ccccc1)C(F)(F)F
Please reason step by step, and put your final answer within `\boxed{}`.

TOOLUSE

Your task is to answer the user's question using available tools. You have access to the following tools:

Name: Axolotl
Description: Collection of axolotl pictures and facts
Documentation:
getRandomAxolotlImage: Retrieve a random axolotl image with information on the image source.
Parameters: {}
Output: Successful response.
- Format: application/json
- Structure: Object{url, source, description}
searchAxolotlImages: Search for axolotl images based on specific criteria such as color, gender, and size.
Parameters: {"color": "string. One of: [wild, leucistic, albino]. The color of the axolotl (e.g., 'wild', 'leucistic', 'albino', etc).", "gender": "string. One of: [male, female]. The gender of the axolotl ('male', 'female').", "size": "string. One of: [small, medium, large]. The size of the axolotl ('small', 'medium', 'large').", "page": "integer. The page number for pagination purposes."}
Output: Successful response.
- Format: application/json
- Structure: Object{results: Array[Object{url, source, description}], pagination: Object{current_page, total_pages, total_results}}
getAxolotlFacts: Retrieve interesting facts about axolotls such as their habits, habitats, and physical characteristics.
Parameters: {"category": "string. One of: [habits, habitat, physical characteristics]. The category of facts to retrieve (e.g., 'habits', 'habitat', 'physical characteristics').", "limit": "integer. The maximum number of facts to return."}
Output: Successful response.
- Format: application/json
- Structure: Array[Object{fact, source}]

Use the following format:
Thought: you should always think about what to do
Action: the action to take, should be one of the tool names.
Action Input: the input to the action, must be in JSON format. All of the action input must be realistic and from the user.

Begin!
Question: Hey, can you show me a random picture of an axolotl?

ZLogic

There are 6 houses, numbered 1 to 6 from left to right, as seen from across the street. Each house is occupied by a different person. Each house has a unique attribute for each of the following characteristics:

- Each person has a unique name: 'Arnold', 'Peter', 'Eric', 'Alice', 'Bob', 'Carol'
- People have unique favorite book genres: 'biography', 'science fiction', 'fantasy', 'mystery', 'romance', 'historical fiction'
- People have unique favorite sports: 'baseball', 'basketball', 'swimming', 'volleyball', 'tennis', 'soccer'
- People own unique car models: 'honda civic', 'ford f150', 'tesla model 3', 'chevrolet silverado', 'bmw 3 series', 'toyota camry'

Clues:
1. Eric is the person who loves mystery books.

```

2. The person who loves tennis is the person who loves fantasy books.
3. The person who loves soccer is directly left of the person who loves science fiction books.
4. There is one house between the person who owns a Honda Civic and the person who loves biography books.
5. Peter is somewhere to the right of Carol.
6. The person who loves tennis is in the first house.
7. The person who owns a Tesla Model 3 is somewhere to the right of the person who loves baseball.
8. Eric is somewhere to the left of the person who loves romance books.
9. The person who owns a Toyota Camry is somewhere to the right of the person who loves romance books.
10. The person who owns a BMW 3 Series is Peter.
11. The person who owns a BMW 3 Series is the person who loves basketball.
12. The person who owns a Tesla Model 3 is directly left of Arnold.
13. Alice and the person who loves volleyball are next to each other.
14. The person who loves historical fiction books is the person who loves soccer.
15. The person who owns a Chevrolet Silverado is not in the first house.
16. The person who loves science fiction books is directly left of the person who loves swimming.

What is Name of the person who lives in House 5?
A. Eric
B. Bob
C. Alice
D. Peter
E. Carol
F. Arnold
Please reason step by step, and put your final answer within \boxed{ }.

```

C.3 MODEL CHAT TEMPLATES

Qwen3-4B-Thinking-2507

```

<|im_start|>user\n{prompt}<|im_end|>
<|im_start|>assistant\n<think>

```

Qwen3-4B-Instruct-2507

```

<|im_start|>user\n{prompt}<|im_end|>
<|im_start|>assistant\n

```

OLMo-3-7B-Instruct

```

<|im_start|>system\nYou are a helpful function-calling AI assistant. You do not currently
↪ have access to any functions. <functions></functions><|im_end|>
<|im_start|>user\n{prompt}<|im_end|>
<|im_start|>assistant\n

```

OLMo-3-7B-Think

```

<|im_start|>system\nYou are OLMo, a helpful function-calling AI assistant built by Ai2.
↪ Your date cutoff is November 2024, and your model weights are available at
↪ https://huggingface.co/allenai. You do not currently have access to any functions.
↪ <functions></functions><|im_end|>
<|im_start|>user\n{prompt}<|im_end|>
<|im_start|>assistant\n<think>

```

C.4 PROMPT TEMPLATES

Prompt for MATH & AIME & Math500 & SCIENCE & GPQA

```

{question}
Please reason step by step, and put your final answer within \boxed{ }.

```

Prompt for TOOLUSE & BFCLv4

```

Your task is to answer the user's question using available tools.
You have access to the following tools:
{tool}

Use the following format:
Thought: you should always think about what to do
Action: the action to take, should be one of the tool names.
Action Input: the input to the action, must be in JSON format. All of the action input must
↪ be realistic and from the user.

Begin!
Question: {question}

```

Table D.1: SDPO performance (Acc@8) against teacher EMA α . No α is globally best.

Trained on	α (%)	AIME	Math500	GPQA	MMLU-R	BFCLv4	LCBv6	Avg
–	–	32.71	73.85	62.94	86.02	76.25	46.73	63.08
MATH	0	56.42	84.29	58.72	85.47	77.71	48.52	68.61
	0.5	57.92	84.52	57.01	84.65	76.46	51.10	68.61
	1	56.04	84.35	58.21	85.31	77.71	51.79	69.06
	2	52.71	84.02	59.53	85.00	78.75	50.35	68.39
	5	34.38	78.97	55.87	83.82	78.12	50.86	63.67
TOOLUSE	0	47.71	78.88	60.73	87.82	76.04	58.02	68.20
	0.5	43.13	77.08	61.49	87.40	72.08	51.68	65.48
	1	41.04	77.18	58.78	87.18	76.04	48.62	64.81
	2	45.00	77.70	58.90	86.75	74.79	53.51	66.11
	5	47.71	78.57	57.39	86.25	77.71	55.55	67.20
SCIENCE	0	55.63	79.25	61.62	87.00	74.79	51.94	68.37
	0.5	53.96	79.10	61.05	87.42	75.21	52.93	68.28
	1	55.84	79.20	62.31	86.92	75.21	51.68	68.53
	2	55.00	80.97	59.47	86.40	76.04	54.92	68.80
	5	53.34	78.42	57.39	84.10	72.29	56.11	66.94

Table D.2: SDPO with and without CoT Performance (Acc@8). CoT enhances TOOLUSE training while impairing MATH and SCIENCE. Larger teacher EMA α generally worsens CoT performance.

α (%)	CoT	Trained on	AIME	Math500	GPQA	MMLU-R	BFCLv4	LCBv6
–	–	–	32.71	73.85	62.94	86.02	76.25	46.73
0	×	MATH	56.42	84.29	58.72	85.47	77.71	48.52
0	✓	MATH	38.13	75.58	58.59	86.22	66.46	56.82
1	×	MATH	56.04	84.89	58.62	85.31	77.71	51.79
1	✓	MATH	37.71	75.58	56.25	86.20	66.88	54.45
5	×	MATH	34.38	78.97	55.87	83.82	78.12	50.86
5	✓	MATH	8.13	69.50	40.53	80.95	34.38	41.76
0	×	SCIENCE	55.63	79.25	61.62	87.00	74.79	51.94
0	✓	SCIENCE	42.08	76.25	55.05	87.08	73.12	57.08
1	×	SCIENCE	55.84	79.20	62.31	86.92	75.21	51.68
1	✓	SCIENCE	37.29	75.60	55.11	86.65	71.67	60.46
5	×	SCIENCE	53.34	78.42	57.39	84.10	72.29	56.11
5	✓	SCIENCE	15.63	71.50	39.96	81.85	60.00	47.77
0	×	TOOLUSE	47.71	78.88	60.73	87.82	76.04	58.48
0	✓	TOOLUSE	53.34	82.28	56.69	84.38	80.62	56.66
1	×	TOOLUSE	41.04	77.18	58.78	87.18	76.04	49.69
1	✓	TOOLUSE	53.55	82.90	56.38	84.55	80.42	59.11
5	×	TOOLUSE	47.71	78.57	57.39	86.25	77.71	55.90
5	✓	TOOLUSE	53.34	83.00	50.06	83.30	77.08	52.50

SDPO Regeneration Prompt

```
{question}
Please reason step by step, and put your final answer within \boxed{}.
Correct solution: {previously correct solution}
Correctly solve the original question.
```

D SUPPLEMENTARY RESULTS

Table D.1 presents the performance scores used to plot Figure 2a, and displays a configuration matrix of EMA teachers $\alpha \in \{0, 0.5\%, 1\%, 2\%, 5\%\}$ and training datasets MATH, TOOLUSE, and SCIENCE. Table D.2 extends Figure 4 for SDPO versus SDPO^{CoT} under more teacher EMA settings and evaluation benchmarks. The large degradation at $\alpha = 1\%, 5\%$ further verifies *confirmation bias*.

Table D.3: Performance (Acc@8) trained on a single domain. Rows indicate training domain and method; columns indicate evaluation benchmarks. GRPO outperforms SDPO at generalization.

	Method	AIME	Math500	SCIENCE	GPQA	TOOLUSE	BFCLv4	ZLogic	MMLU-R
Olmo-3-7B-Instruct-DPO									
	–	19.17	61.30	32.50	30.24	42.46	76.46	60.48	69.92
MATH	GRPO	18.96	64.82	35.24	35.23	40.62	71.88	64.95	72.00
	SDPO ₀	13.75	51.24	21.93	24.94	34.38	57.08	34.31	58.96
Olmo-3-7B-Think-DPO									
	–	46.04	71.53	30.89	36.30	9.19	3.96	59.18	79.35
MATH	GRPO	51.67	74.95	37.56	44.32	9.38	3.33	61.58	79.47
	SDPO _{5%}	41.67	72.90	28.15	33.71	8.09	0.00	43.88	72.20
	SDPO ₀	51.25	75.95	32.98	38.45	8.46	0.62	54.40	79.00
	SDPO ₀ ^{COT}	54.59	74.52	41.79	43.81	7.17	0.00	61.78	78.53

Table D.4: Representative gain, retention, and forgetting cases. Intermediate-distance targets exhibit the clearest forgetting, while source-like and very distant targets are often stable or improve.

	Train	Test	MMD·10 ⁻²	$\Delta_{SDPO_{5\%}}$	Δ_{SDPO_0}	Gain	Retention	Forgetting
Near	MATH	AIME	0.64	↑ 1.67	↑ 23.71	✓		
	MATH	Math500	1.38	↑ 5.12	↑ 10.44	✓		
	CODING	LCBv6	3.74	↑ 11.74	↑ 13.39	✓		
	TOOLUSE	BFCLv4	4.75	↑ 1.46	↓ 0.21		✓	
Intermediate	MATH	GPQA	6.74	↓ 7.07	↓ 4.22			✓
	SCIENCE	GPQA	11.42	↓ 5.55	↓ 1.32			✓
	TOOLUSE	GPQA	18.89	↓ 5.55	↓ 2.21			✓
	CODING	GPQA	11.90	↓ 4.61	↓ 2.65			✓
	MATH	ZLogic	21.98	↓ 7.94	↓ 7.12			✓
	SCIENCE	BFCLv4	25.16	↓ 3.96	↓ 1.46			✓
Distant	SCIENCE	ZLogic	32.66	↑ 9.06	↑ 8.18	✓		
	TOOLUSE	ZLogic	43.33	↑ 2.40	↑ 4.86	✓		
	CODING	ZLogic	31.66	↓ 1.09	↑ 0.96		✓	

Table D.5: SDPO_{5%} performance (Acc@8) continually trained with various data at 2nd & 3rd stage. Altering order only partially mitigates collapse; SDPO still degrades significantly under data shifts.

	Trained on	MATH		SCIENCE		TOOLUSE		CODING		LOGIC	KNOW.
		AIME	Math500	SCIENCE	GPQA	TOOLUSE	BFCLv4	CODING	LCBv6	ZLogic	MMLU-R
	–	32.71	73.85	50.54	62.94	48.90	76.25	56.39	46.73	81.62	86.02
SDPO _{5%}	1. MATH	34.38	78.97	45.00	55.87	53.31	78.12	52.37	50.86	73.68	83.82
	2. SCIENCE	21.25	73.95	76.37	47.98	30.15	51.25	59.33	48.82	61.88	82.12
	3. TOOLUSE	0.00	0.00	0.18	0.13	0.00	0.00	0.00	0.00	0.10	0.12
	2. SCIENCE	21.25	73.95	76.37	47.98	30.15	51.25	59.33	48.82	61.88	82.12
	3. CODING	17.09	69.53	74.52	42.93	0.74	0.00	54.79	54.01	58.00	80.50
	2. TOOLUSE	22.29	73.82	38.63	42.36	12.68	9.79	–	–	56.95	77.35
	3. SCIENCE	10.42	66.50	72.56	36.55	13.60	20.62	–	–	49.20	77.75

For single-domain training, Table D.3 examines SDPO and GRPO for Olmo-3 series. For continual training, Table D.6 presents detailed scores for Figure 6 and table 3 across each stage, with a data-ordering ablation on collapsed SDPO_{5%} in Table D.5. We supplement diagnostics (Figure 7) with additional matrices (Figure D.1), more training (Figure D.2), and more benchmarks (Figure D.3).

Table D.6: SDPO versus GRPO performance (Acc@8) continually trained. $SDPO_0$ and $SDPO_{5\%}$ denote SDPO with teacher EMA 0 and 5%, respectively. $SDPO_0^{CoT}$ denotes $SDPO_0$ with CoT distilled. $SDPO_{5\%}^{M25}$ denotes $SDPO_{5\%}$ randomly masked out 25% tokens. S- $SDPO_0$ denotes SDPO with restart-and-freeze teacher strategy ($\gamma = 1$).

Trained on	MATH		SCIENCE		TOOLUSE		CODING		LOGIC	KNOW.	Avg	
	AIME	Math500	SCIENCE	GPQA	TOOLUSE	BFCLv4	CODING	LCBv6	ZLogic	MMLU-R		
Qwen3-4B-Thinking-2507												
–	32.71	73.85	50.54	62.94	48.90	76.25	56.39	46.73	81.62	86.02	61.60	
GRPO	1. MATH	44.67	76.84	51.33	62.83	50.48	78.33	60.73	52.93	80.82	85.74	64.47
	2. SCIENCE	54.88	78.28	56.39	64.23	52.06	76.67	60.06	50.17	88.23	87.91	66.89
	3. TOOLUSE	56.21	79.49	60.64	63.99	58.71	80.83	60.49	52.23	90.10	88.54	69.12
	4. CODING	57.29	79.88	61.85	64.58	57.72	81.25	64.68	53.56	89.52	88.70	69.90
$SDPO_0$	1. MATH	56.42	84.29	47.36	58.72	52.06	77.71	57.58	48.52	74.50	85.47	64.26
	2. SCIENCE	49.29	81.75	65.70	59.03	48.49	76.67	58.35	56.88	87.17	86.38	66.97
	3. TOOLUSE	41.13	80.02	70.37	54.71	37.35	30.83	56.39	54.25	80.03	84.92	59.00
	4. CODING	34.38	77.40	70.71	51.83	9.93	11.25	61.77	51.89	72.15	84.20	52.55
$SDPO_0^{CoT}$	1. MATH	38.13	75.58	49.88	58.59	48.35	66.46	52.70	56.82	77.00	86.22	60.97
	2. SCIENCE	27.71	74.65	70.89	54.10	42.83	61.67	47.28	46.28	76.65	85.95	58.80
	3. TOOLUSE	21.04	71.88	70.71	49.37	38.60	36.04	41.94	44.88	66.22	84.38	52.51
	4. CODING	14.17	69.53	68.69	43.24	23.35	13.54	30.01	24.32	59.95	81.62	42.84
$SDPO_{5\%}$	1. MATH	34.38	78.97	45.00	55.87	53.31	78.12	52.37	50.86	73.68	83.82	60.64
	2. SCIENCE	21.25	73.95	76.37	47.98	30.15	51.25	59.33	48.82	61.88	82.12	55.31
	3. TOOLUSE	0.00	0.00	0.18	0.13	0.00	0.00	0.00	0.00	0.10	0.12	0.05
	4. CODING	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.02	0.18	0.02
$SDPO_{5\%}^{M25}$	1. MATH	37.09	79.88	42.32	56.12	52.39	74.58	54.71	51.51	71.80	84.35	60.48
	2. SCIENCE	18.54	72.38	72.32	45.08	23.16	23.96	57.73	48.75	54.12	81.90	49.79
	3. TOOLUSE	0.00	0.00	5.42	1.96	0.00	0.00	0.00	0.00	1.50	9.65	1.85
$SDPO_{5\%}^{M50}$	1. MATH	39.17	78.53	43.99	55.81	50.00	76.67	54.01	51.75	75.20	84.72	60.98
	2. SCIENCE	25.84	76.10	77.98	51.52	35.48	44.58	48.18	49.56	72.05	82.98	56.43
	3. TOOLUSE	0.00	0.12	34.17	18.62	0.00	0.00	0.09	0.00	55.47	43.38	15.18
$SDPO_{5\%}^{M75}$	1. MATH	37.50	79.05	43.15	56.12	53.68	80.00	53.53	51.93	75.22	84.85	61.50
	2. SCIENCE	18.55	73.05	70.95	44.00	42.83	63.75	59.02	50.39	56.53	81.73	55.31
	3. TOOLUSE	0.63	28.48	54.70	28.22	0.00	0.00	8.18	7.31	45.55	62.02	23.51
Qwen3-4B-Instruct-2507												
–	50.00	77.10	24.23	14.71	59.01	83.12	–	–	48.00	56.75	51.61	
GRPO	1. MATH	56.00	79.94	25.94	17.23	58.97	83.33	–	–	49.40	57.60	53.55
	2. SCIENCE	56.88	81.12	59.35	58.84	57.72	82.92	–	–	84.22	86.52	70.95
	3. TOOLUSE	58.34	81.88	60.54	59.91	65.99	75.42	–	–	84.80	87.85	71.84
$SDPO_0$	1. MATH	35.84	71.15	12.20	5.49	59.74	82.29	–	–	9.15	2.97	34.85
	2. SCIENCE	27.08	64.53	11.79	8.90	57.90	82.92	–	–	11.70	9.48	34.29
	3. TOOLUSE	20.63	54.25	7.14	4.17	0.00	0.00	–	–	1.48	5.35	11.63
SFT	1. MATH	0.42	20.78	27.38	23.99	52.94	78.12	–	–	2.58	66.55	34.09
	2. SCIENCE	0.00	0.00	72.74	30.74	0.00	7.92	–	–	30.22	69.30	26.36
	3. TOOLUSE	0.00	0.00	0.00	0.00	0.00	0.00	–	–	0.00	0.00	0.00

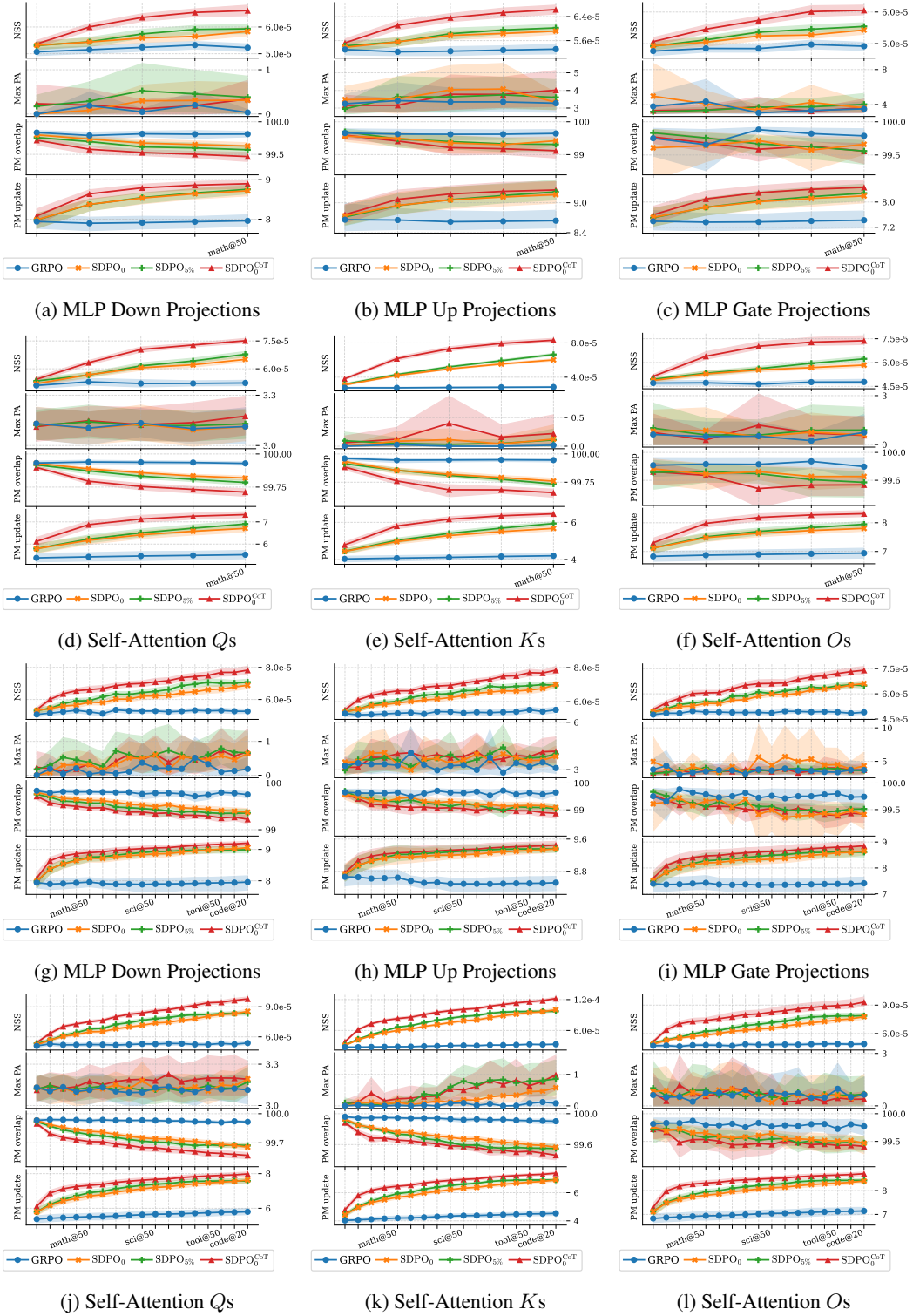


Figure D.1: Parameter drift (y -axis) against training steps (x -axis), with metrics averaged across all layers with a 95% confidence interval. (a-f) MATH training. (g-l) continual training.

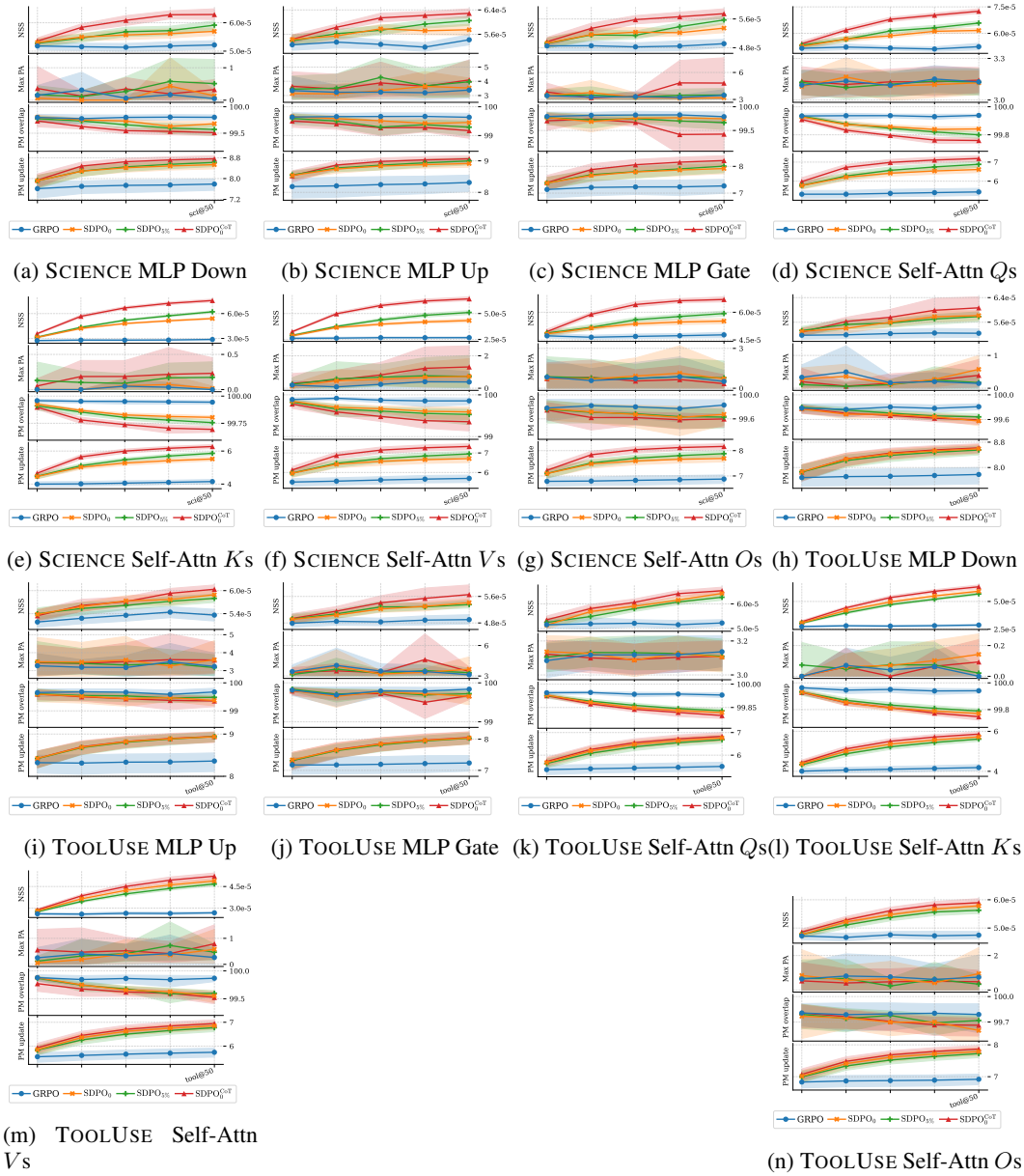


Figure D.2: Parameter drift (y -axis) against SCIENCE or TOOLUSE training steps (x -axis, $t = 10, 20, \dots, 50$), with metrics averaged across all layers with a 95% confidence interval.

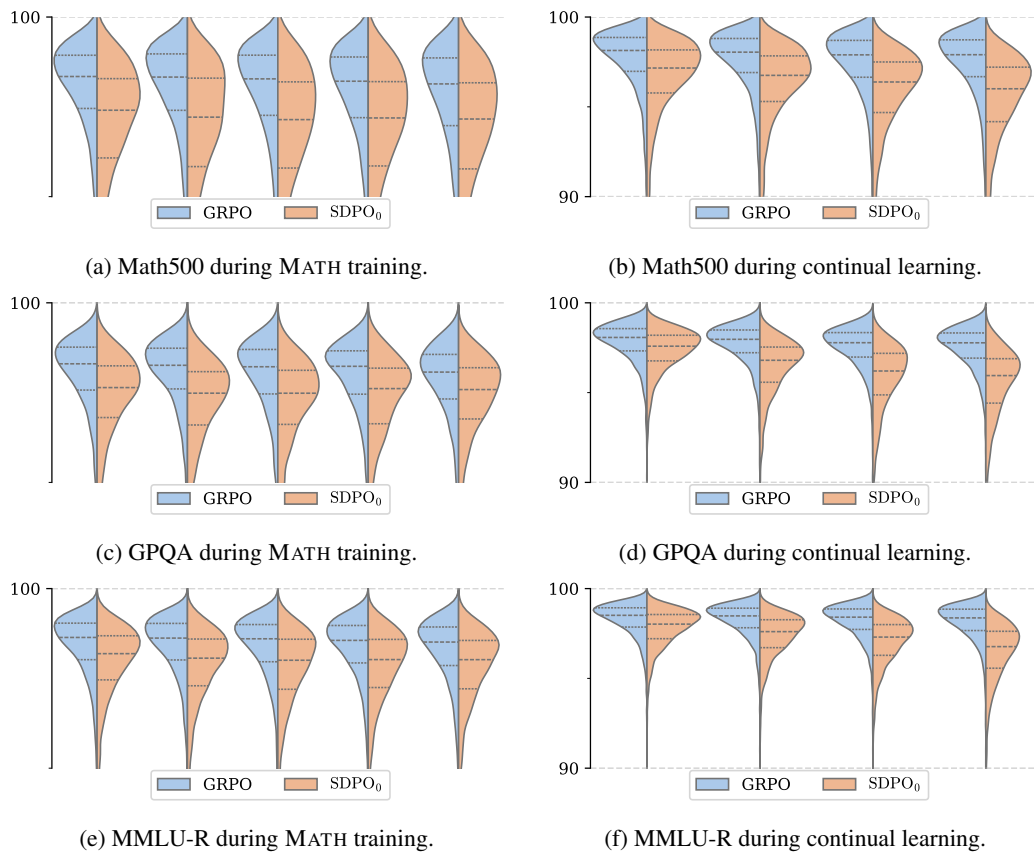


Figure D.3: Cosine similarity (%) of rollout(@8) embeddings before and after GRPO (in blue) versus SDPO₀ (in orange) training. For MATH training, x -axis denote checkpoints at steps $k = 10, 20, \dots, 50$ from left to right; for continual learning, x -axis denote checkpoints after tasks MATH, SCIENCE, TOOLUSE, and CODING from left to right.

E GENERATION SAMPLES

AIME Qwen3-4B-Thinking-2507 Generation

This is a complex or challenging question, and it is difficult to provide a direct and correct answer. I need to think about it. Well, let's tackle this problem step by step. So, first, let's make sure I understand what's given and what's being asked. Aya does a 9 km walk every morning, but wait, hold on. Wait, no, wait: "a 9-kilometer-long walk" does that mean the walk itself is 9 km, not including the coffee shop? Wait, the problem says "the walk takes her 4 hours, including t minutes spent in the coffee shop." Hmm, okay, so the "walk" here is the time from when she starts walking until she gets to the coffee shop? Wait, no, maybe better to parse the problem carefully.

"Every morning Aya goes for a 9-kilometer-long walk and stops at a coffee shop afterwards." So, the walk is 9 km long, then she stops at a coffee shop (so the coffee shop is the endpoint of the walk? Maybe the walk is the 9 km journey to the coffee shop, and then she spends t minutes there? Wait, but the problem says "the walk takes her 4 hours, including t minutes spent in the coffee shop." Wait, that's confusing. If the walk is the 9 km, then why is the time including time spent in the coffee shop? Maybe the problem is using "walk" to mean the entire activity? Wait, no, let's check the problem statement again:

"When she walks at a constant speed of s kilometers per hour, the walk takes her 4 hours, including t minutes spent in the coffee shop." Oh! Wait, maybe "the walk" here is a misnomer, but actually, the total time from when she starts walking until she finishes at the coffee shop (i.e., the time to walk the 9 km plus the time she spends at the coffee shop) is 4 hours? Wait, but the problem says "the walk takes her 4 hours, including t minutes spent in the coffee shop." Hmm, maybe the problem is translated or phrased a bit awkwardly, but let's think in terms of variables.

Let's define:

Let's let the time she spends walking (i.e., the time actually moving) be W hours, and the time spent at the coffee shop be T hours? Wait, no, the problem says t minutes, so maybe better to use minutes for the coffee shop time to avoid decimals, but maybe keep it in hours first.

Wait, let's try to set up equations based on the two scenarios.

First scenario: speed = s km/h, total time (walk + coffee shop) = 4 hours, which includes t minutes at coffee shop. Wait, but hold on. If the walk is 9 km long, then the walking time should be distance over speed, right? So maybe the "walk" here is the 9 km, so the time she spends walking (moving) is $9/s$ hours, and then she spends t minutes at the coffee shop, so total time from start to finish (including coffee shop) is $(9/s) + (t/60)$ hours. The problem says "the walk takes her 4 hours, including t minutes spent in the coffee shop." Wait, maybe the problem is using "walk" to mean the entire trip? Like, "the walk" is the activity from when she starts walking until she stops at the coffee shop, so the duration of that activity is 4 hours, which consists of the time walking (9 km) plus time at coffee shop (t minutes). That might be the right interpretation.

Yes, that makes sense. So "the walk takes her 4 hours" the entire walk activity (which includes walking the 9 km and then stopping at coffee shop) takes 4 hours total, with t minutes of that time spent at the coffee shop. So the rest of the time (4 hours minus t minutes) is spent actually walking the 9 km.

Similarly, second scenario: speed = s + 2 km/h, total walk activity time = 2 hours 24 minutes, which includes t minutes at coffee shop, so walking time is $2\text{h}24\text{m} - t$ minutes.

Okay, that seems to be the key interpretation. Let's formalize this.

Let me convert all times to hours to make the equations easier, or maybe keep t in minutes and convert the total times to minutes? Let's see. Let's try hours first, then convert t to hours.

First scenario:

Total time (including coffee shop) = 4 hours = 240 minutes.

Time spent at coffee shop = t minutes = $t/60$ hours.

Therefore, time spent walking = total time - coffee shop time = $240 - t$ minutes = $(240 - t)/60$ hours.

But she walks 9 km at speed s km/h, so walking time = $9 / s$ hours.

Therefore, we have equation 1: $9 / s = (240 - t)/60$.

Second scenario:

Total time (including coffee shop) = 2 hours 24 minutes = $2 + 24/60 = 2.4$ hours = 144 minutes.

Time spent at coffee shop = still t minutes (the problem says "including t minutes spent in the coffee shop" in both cases, so t is the same in both scenarios she spends t minutes at coffee shop regardless of walking speed, which is good, that's a constant we need to find or relate).

Therefore, time spent walking = total time - coffee shop time = $144 - t$ minutes = $(144 - t)/60$ hours.

She walks 9 km at speed s + 2 km/h, so walking time = $9 / (s + 2)$ hours.

Therefore, equation 2: $9 / (s + 2) = (144 - t)/60$.

Okay, great! So now we have two equations with two variables: s (km/h) and t (minutes). We need to solve for s and t, then find the total time when speed is s + 1/2 km/h, which would be (walking time) + t minutes, so total time in minutes is $(9 / (s + 1/2) * 60) + t$, right? Wait, let's confirm: walking time at speed s + 1/2 is $9 / (s + 1/2)$ hours, which is $60 * (9 / (s + 1/2))$ minutes, then add t minutes for coffee shop, so total time in minutes is $540 / (s + 1/2) + t$. Maybe better to keep everything in minutes to avoid confusion later.

Let's write equations 1 and 2 again, but let's multiply both sides by 60 to eliminate denominators, that
 \hookrightarrow might make it cleaner.

Equation 1: $9 / s = (240 - t) / 60$ multiply both sides by 60s: $540 = s(240 - t)$ $540 = 240s -$
 \hookrightarrow st . Let's call this Eq1.

Equation 2: $9 / (s + 2) = (144 - t) / 60$ multiply both sides by $60(s + 2)$: $540 = (s + 2)(144 - t)$.
 \hookrightarrow Let's call this Eq2.

Okay, now we have two linear equations in terms of s and t ? Wait, no, actually, Eq1 and Eq2 are both
 \hookrightarrow linear in s and t , yes! Let's check:

Eq1: $240s - st = 540$ $s(240 - t) = 540$

Eq2: $(s + 2)(144 - t) = 540$

Oh, nice! Both equal 540, so we can set them equal to each other to eliminate the 540.

So, $s(240 - t) = (s + 2)(144 - t)$

Let's expand the right-hand side: $s(144 - t) + 2(144 - t) = 144s - st + 288 - 2t$

Left-hand side: $240s - st$

So set equal: $240s - st = 144s - st + 288 - 2t$

Wait, look, we can subtract $(-st)$ from both sides? Wait, let's subtract $144s - st$ from both sides to
 \hookrightarrow simplify.

Left side minus $(144s - st) = (240s - st) - 144s + st = 96s$

Right side minus $(144s - st) = (144s - st + 288 - 2t) - 144s + st = 288 - 2t$

So we get $96s = 288 - 2t$

Oh, that's a nice simplification! The st terms canceled out, which is good, we didn't want to deal with
 \hookrightarrow quadratic equations. Let's write that as:

$96s = 288 - 2t$ let's divide both sides by 2 to simplify: $48s = 144 - t$ $t = 144 - 48s$. Let's
 \hookrightarrow call this Eq3.

Okay, great, now we can plug Eq3 into one of the original equations to solve for s . Let's pick Eq1: $s(240$
 \hookrightarrow $- t) = 540$.

First, compute $240 - t$ using Eq3: $240 - (144 - 48s) = 240 - 144 + 48s = 96 + 48s$.

So Eq1 becomes: $s(96 + 48s) = 540$.

Wait, hold on: is that right? Wait, Eq1 was $s(240 - t) = 540$, yes, and we found $240 - t = 96 + 48s$, so
 \hookrightarrow yes, $s(96 + 48s) = 540$.

Let's write that as $48s^2 + 96s - 540 = 0$. Maybe we can simplify this equation by dividing all terms
 \hookrightarrow by 24? Let's see: $48/24=2$, $96/24=4$, $540/24=22.5$, which is not integer. Divide by 12: $4s^2 + 8s -$
 \hookrightarrow $45 = 0$. Wait, $48/12=4$, $96/12=8$, $540/12=45$, yes! Perfect, that's integer coefficients.

So $4s^2 + 8s - 45 = 0$. Now let's solve for s using quadratic formula.

Quadratic in standard form: $as^2 + bs + c = 0$, solutions $s = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$

Here, $a=4$, $b=8$, $c=-45$.

Discriminant $D = b^2 - 4ac = 64 - 4 \cdot 4 \cdot (-45) = 64 + 720 = 784$. Oh, nice, 784 is a perfect square!
 \hookrightarrow $\sqrt{784} = 28$.

So $s = \frac{-8 \pm 28}{8}$. Now, speed can't be negative, so we take the positive root: $(-8 + 28) / 8 = 20 / 8$
 \hookrightarrow $= 5 / 2 = 2.5$ km/h? Wait, hold on, wait: 20 divided by 8 is 2.5, yes, that's 2.5 km/h. Wait, but let's
 \hookrightarrow check if that makes sense with the times.

Wait, but hold on, let's make sure we didn't make a mistake in setting up the equations, because if s is
 \hookrightarrow 2.5 km/h, then walking time for 9 km is $9 / 2.5 = 3.6$ hours = 216 minutes, so total time including
 \hookrightarrow coffee shop is $216 + t$ minutes = 4 hours = 240 minutes, so $t = 24$ minutes. Let's check the second
 \hookrightarrow scenario: speed $s + 2 = 4.5$ km/h, walking time = $9 / 4.5 = 2$ hours = 120 minutes, total time including
 \hookrightarrow coffee shop is $120 + t = 120 + 24 = 144$ minutes = 2 hours 24 minutes, which matches the problem
 \hookrightarrow statement! Oh, that's perfect! So that checks out.

Wait, so why did I doubt? Let's go through the algebra again to make sure I didn't skip a step wrong.

We had Eq1: $s(240 - t) = 540$ (from $9/s = (240 - t)/60$ $9 \cdot 60 = s(240 - t)$ $540 = s(240 - t)$,
 \hookrightarrow yes, that's correct, $9 \cdot 60$ is 540, right).

Eq2: $(s + 2)(144 - t) = 540$, same way: $9/(s+2) = (144 - t)/60$ $9 \cdot 60 = (s+2)(144 - t)$ $540 =$
 \hookrightarrow $(s+2)(144 - t)$, correct.

Set equal: $s(240 - t) = (s + 2)(144 - t)$. Expanded RHS: $s(144 - t) + 2(144 - t) = 144s - st + 288 - 2t$.
 \hookrightarrow LHS: $240s - st$. Subtract $144s - st$ from both sides: $96s = 288 - 2t$ $48s = 144 - t$ $t = 144$
 \hookrightarrow $- 48s$, correct.

Then plug into Eq1: $s(240 - (144 - 48s)) = s(96 + 48s) = 540$ $48s^2 + 96s - 540 = 0$
 \hookrightarrow divide by 12: $4s^2 + 8s - 45 = 0$, correct.

Quadratic formula: $s = \frac{-8 \pm \sqrt{64 + 720}}{8} = \frac{-8 \pm \sqrt{784}}{8} = \frac{-8 \pm 28}{8}$.
 \hookrightarrow Positive solution: $20/8 = 2.5$, yes, 2.5 km/h. Then $t = 144 - 48s = 144 - 48 \cdot (2.5) = 144 - 120 = 24$
 \hookrightarrow minutes. Perfect, matches the check we did earlier.

Wait, so in the check, when $s=2.5$, walking time is $9/2.5=3.6$ hours=216 minutes, total time=216 +24=240
 \hookrightarrow minutes=4 hours, correct. Second scenario: $s+2=4.5$, walking time= $9/4.5=2$ hours=120 minutes, total
 \hookrightarrow time=120+24=144 minutes=2h24m, correct. So all checks out.

Now, the problem asks: Suppose Aya walks at $s + 1/2$ kilometers per hour. Find the number of minutes the
 \hookrightarrow walk takes her, including the t minutes spent in the coffee shop.

So first, $s + 1/2 = 2.5 + 0.5 = 3$ km/h. Wait, that's nice! 2.5 is $5/2$, so $5/2 + 1/2 = 6/2 = 3$ km/h. Good,
 \hookrightarrow so speed is 3 km/h.

Now, walking time at 3 km/h for 9 km: $9 / 3 = 3$ hours. Then add t minutes for coffee shop, which we found
 \hookrightarrow $t=24$ minutes.

Wait, hold on! Wait, the problem says "the number of minutes the walk takes her, including the t minutes
 \hookrightarrow spent in the coffee shop." So total time in minutes is (walking time in hours * 60) + t .

Walking time at 3 km/h: $9 \text{ km} / 3 \text{ km/h} = 3$ hours = 180 minutes. Then add $t=24$ minutes, so total time=180
 \hookrightarrow +24=204 minutes.

Wait, is that it? Let me make sure we didn't misinterpret "the walk takes her...". Wait, in the problem
 \hookrightarrow statement, when she walks at s km/h, the walk takes 4 hours including t minutes at coffee shop. So in
 \hookrightarrow our interpretation, "the walk" (the activity) takes total time = walking time + coffee time. So yes,
 \hookrightarrow for the new speed, it's walking time (at $s+1/2$) + t minutes, which is what we did.

But let's just go through the problem statement again to be 100% sure we didn't flip something.

"Every morning Aya goes for a 9-kilometer-long walk and stops at a coffee shop afterwards." So the walk is
 \hookrightarrow 9 km, then she stops at coffee shop (so coffee shop is after the walk? Wait, maybe that's the
 \hookrightarrow confusion! Maybe "the walk" is 9 km, then she goes to coffee shop, so the time to walk the 9 km is
 \hookrightarrow separate from the time at coffee shop, and the problem says "the walk takes her 4 hours, including t
 \hookrightarrow minutes spent in the coffee shop" wait, that would be odd, because if the walk is the 9 km, it
 \hookrightarrow shouldn't include coffee shop time. Wait, but in our check, it worked out perfectly, so maybe the
 \hookrightarrow problem statement is a bit poorly worded, but the math checks out with the interpretation that total
 \hookrightarrow time from start of walk (when she begins walking) until she finishes at coffee shop is 4 hours (which
 \hookrightarrow includes the time she's walking the 9 km plus time at coffee shop).

Wait, another way: suppose "the walk" refers to the 9 km journey, so time to walk 9 km is some time, then
 \hookrightarrow she spends t minutes at coffee shop, so total time from when she starts walking until she leaves
 \hookrightarrow coffee shop is (walking time) + t minutes. The problem says "the walk takes her 4 hours, including t
 \hookrightarrow minutes spent in the coffee shop" maybe "takes her" is the duration of the entire activity
 \hookrightarrow (walking + coffee), so yes, that's how we took it, and the check worked, so probably that's the right
 \hookrightarrow interpretation.

Let's see, if we had misinterpreted and thought that the coffee shop time is part of the walk time, but
 \hookrightarrow no, the problem says "including t minutes spent in the coffee shop", so the t minutes are added to the
 \hookrightarrow walk time to get the total time.

Wait, in our check, with $s=2.5$ km/h, walking time (just walking) is 3.6 hours=216 minutes, total time with
 \hookrightarrow coffee is 240 minutes, so $t=24$ minutes. Second scenario: walking time 2 hours=120 minutes, total time
 \hookrightarrow 144 minutes, so $t=24$ minutes, same t , which the problem states, so that's consistent.

So when speed is 3 km/h, walking time is $9/3=3$ hours=180 minutes, total time=180 +24=204 minutes.

Wait, but let's try to do it without plugging in numbers early, to verify.

We had $s=2.5$, $t=24$, so $s + 1/2=3$, walking time= $9/3=3$ hours=180 min, total=180+24=204.

Alternatively, let's see if we can find the total time without finding s and t individually, maybe as a
 \hookrightarrow check.

Let s denote for a given speed v (km/h), walking time is $9/v$ hours, total time $T(v) = 9/v + t/60$
 \hookrightarrow hours (since t minutes is $t/60$ hours).

In first case, $v=s$, $T(s)=4$ hours $9/s + t/60 = 4$ $9/s = 4 - t/60$ which is what we had
 \hookrightarrow (since $4 - t/60 = (240 - t)/60$ hours).

Second case, $v=s+2$, $T(s+2)=2 + 24/60=2.4$ hours $9/(s+2) + t/60=2.4$ $9/(s+2)=2.4 - t/60=(144 -$
 \hookrightarrow $t)/60$, same as before.

So we can write the two equations as:

$$9/s = 4 - t/60 \dots(1a)$$

$$9/(s+2) = 2.4 - t/60 \dots(2a)$$

Let s let $u = t/60$, so u is time in hours at coffee shop. Then equations become:

$$9/s = 4 - u \dots(1b)$$

$$9/(s+2) = 2.4 - u \dots(2b)$$

Now, subtract equation (2b) from equation (1b):

$$9/s - 9/(s+2) = (4 - u) - (2.4 - u) = 1.6$$

Simplify left side: $9[(s+2 - s)/(s(s+2))] = 9[2/(s(s+2))] = 18/(s(s+2))$

So $18/(s(s+2)) = 1.6 = 8/5$ (converting 1.6 to fraction: $8/5$)

Therefore, $s(s+2) = 18 * (5/8) = 90/8 = 45/4 = 11.25$

Wait, $s(s+2)=45/4$ $s^2 + 2s - 45/4=0$ multiply by 4: $4s^2 + 8s - 45=0$, which is exactly
 \hookrightarrow the quadratic we had earlier! Good, so that's a good consistency check.

Then from (1b): $u=4 - 9/s$. We know $s=2.5=5/2$, so $9/s=9/(5/2)=18/5=3.6$, so $u=4 - 3.6=0.4$ hours=24 minutes,
 \hookrightarrow correct.

Now, for speed $v=s + 1/2=5/2 + 1/2=3$ km/h, total time $T(v)=9/v + u=9/3 + 0.4=3 + 0.4=3.4$ hours. Convert to
 \hookrightarrow minutes: $3.4 \cdot 60=204$ minutes. Perfect, same answer.

So that's reassuring. The key was setting up the equations correctly, which we verified with two methods
 \hookrightarrow and a check, so the total time is 204 minutes.

Wait, just to make sure there isn't a trick in the problem statement: "the walk takes her 4 hours,
 \hookrightarrow including t minutes spent in the coffee shop". Does "the walk" refer to the coffee shop trip? No, the
 \hookrightarrow coffee shop is where she stops after the walk, so the walk is the 9 km, then coffee shop is separate.
 \hookrightarrow But the problem says the time for the walk (activity) includes coffee shop time, which is a bit
 \hookrightarrow confusing, but since the math checks out with the problem's given times (4 hours total with t minutes,
 \hookrightarrow 2h24m total with t minutes), and we found consistent t , it's probably the right interpretation.

Another possible misinterpretation: maybe t minutes is the time she spends walking? But no, the problem
 \hookrightarrow says " t minutes spent in the coffee shop", so t is coffee time.

Wait, let's suppose for a second that we misread and thought the 9 km is the time or something, but no, it
 \hookrightarrow says "9-kilometer-long walk", so distance is 9 km.

Wait, let's see if the problem says "the walk takes her 4 hours" if "the walk" is the 9 km, then
 \hookrightarrow walking time is 4 hours, but then she spends t minutes at coffee shop, so total time is $4h + t$ min,
 \hookrightarrow but the problem says "the walk takes her 4 hours, including t minutes...", so "including" implies that
 \hookrightarrow the t minutes are part of the 4 hours, so the walk's duration is 4 hours, which comprises walking and
 \hookrightarrow coffee time. So yes, our original interpretation is correct.

I think the key was not overcomplicating it and setting up the equations with total time = walking time +
 \hookrightarrow coffee time, which gave us two equations that simplified nicely, and the check confirmed the values of
 \hookrightarrow s and t , so we didn't have to worry about extraneous solutions or wrong interpretations.

So to recap:

- Found $s=2.5$ km/h, $t=24$ minutes via equations, verified with time calculations.
- New speed: $2.5 + 0.5=3$ km/h.
- Walking time for 9 km at 3 km/h: 3 hours=180 minutes.
- Total time including coffee: $180 + 24=204$ minutes.

All steps check out, multiple methods confirm the answer. I think that's solid.

I think sometimes when problems say "the walk takes X time including Y time", it's easy to get confused
 \hookrightarrow about what's included, but by testing with the numbers we got (which matched the problem's time
 \hookrightarrow descriptions), we could be confident in the setup.

Just to write the final answer clearly: 204 minutes.

****Final Answer**** $\boxed{204}$ $\nless than /think>$

We are given that Aya walks a ****9-kilometer-long walk**** and stops at a coffee shop after. The problem
 \hookrightarrow states:

- When she walks at a constant speed of s km/h, the ****total time**** (including t minutes spent in
 \hookrightarrow the coffee shop) is ****4 hours****.
 \hookrightarrow - When she walks at $s + 2$ km/h, the ****total time**** (including s
 \hookrightarrow t minutes at the coffee shop) is ****2 hours and 24 minutes****.

We are to find how many ****minutes**** the walk takes her when she walks at $s + \frac{1}{2}$ km/h,
 \hookrightarrow including the t minutes spent at the coffee shop.

Step 1: Convert all times into minutes for consistency

Let's define:
 \hookrightarrow - Total time in minutes for first scenario: $4 \times 60 = 240$ minutes
 \hookrightarrow time in minutes for second scenario: $2 \times 60 + 24 = 144$ minutes

Let t be the number of minutes spent at the coffee shop (same in both scenarios)
 \hookrightarrow Then the ****walking time**** in each case is:
 \hookrightarrow - First scenario: $240 - t$ minutes
 \hookrightarrow - Second scenario: $144 - t$ minutes

Since she walks 9 km in each case, we can write equations based on speed and time:

$$s = \frac{\text{Distance}}{\text{Time}} \implies \text{Time} = \frac{\text{Distance}}{s}$$

So for the first scenario:

$$\frac{9}{s} = \frac{240 - t}{60} \implies 540 = s(240 - t)$$

For the second scenario:

$$\frac{9}{s + 2} = \frac{144 - t}{60} \implies 540 = (s + 2)(144 - t)$$

Step 2: Solve the system of equations

From (1): $s(240 - t) = 540$

From (2): $(s + 2)(144 - t) = 540$

```

Set them equal:

$$s(240 - t) = (s + 2)(144 - t)$$

Expand both sides:
Left:  $240s - st$ 
Right:  $s(144 - t) + 2(144 - t) = 144s - st + 288 - 2t$ 
Subtract  $144s - st$  from both sides:

$$288 - 2t = 144 - t \quad \rightarrow \quad t = 144 - 48s$$

Now plug (3) into equation (1):

$$s(240 - (144 - 48s)) = 540 \quad \rightarrow \quad s(96 + 48s) = 540 \quad \rightarrow \quad 48s^2 + 96s - 540 = 0$$

Divide by 12:

$$4s^2 + 8s - 45 = 0$$

Use quadratic formula:

$$s = \frac{-8 \pm \sqrt{64 + 720}}{8} = \frac{-8 \pm \sqrt{784}}{8} = \frac{-8 \pm 28}{8}$$

Take positive root:

$$s = \frac{20}{8} = 2.5 \text{ km/h}$$

Now compute  $t$  using (3):

$$t = 144 - 48 \cdot 2.5 = 144 - 120 = 24 \text{ minutes}$$

---
### Step 3: Compute total time at new speed
New speed:  $s + \frac{1}{2} = 2.5 + 0.5 = 3 \text{ km/h}$ 
Walking time for 9 km:

$$\frac{9}{3} = 3 \text{ hours} = 180 \text{ minutes}$$

Add coffee shop time  $t = 24$  minutes:

$$\text{Total time} = 180 + 24 = \boxed{204} \text{ minutes}$$

---
### Final Answer

$$\boxed{204}$$


```

TOOLUSE Qwen3-4B-Thinking-2507 Generation

```

Okay, the user is asking for a random picture of an axolotl. Let me check the available tools.

First, there's the Axolotl tool with two main functions: getRandomAxolotlImage and searchAxolotlImages.
↪ Since the user wants a random picture, the getRandomAxolotlImage tool seems perfect here. The search
↪ tool requires parameters like color, gender, size, which the user didn't specify, so that's not
↪ needed.

The getRandomAxolotlImage has no parameters, so the Action Input should be an empty JSON object. I need to
↪ make sure the response is in the correct JSON format as per the documentation. The output should be an
↪ object with url, source, and description.

I should call getRandomAxolotlImage with no parameters. Let me structure the Action correctly. The user
↪ didn't mention any specific criteria, so no need for the search tool. Just use the random one.
</think>

Thought: The user is asking for a random picture of an axolotl. The appropriate tool for this is
↪ getRandomAxolotlImage, which requires no parameters.

Action: getRandomAxolotlImage
Action Input: {}

```