
AsyncOPD: How Stale Can On-Policy Distillation Be?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 On-policy distillation (OPD) is an increasingly important way to post-train large
2 language models (LLMs), but, like reinforcement learning (RL), it relies on rollouts
3 from the model being optimized. For reasoning workloads, rollout generation can
4 dominate training time. Asynchronous RL alleviates this bottleneck by decoupling
5 rollout generation from learner updates, but doing so introduces stale-policy data;
6 prior work studies how to stabilize learning from such data. However, it remains un-
7 derexplored whether these asynchronous RL ideas and stale-data solutions transfer
8 to OPD, and what OPD-specific constraints arise. To address this gap, we provide
9 the first systematic study of staleness in asynchronous OPD. We first show that KL
10 direction changes the stale-data problem: teacher-weighted forward KL is robust to
11 stale rollouts, whereas student-weighted reverse KL is vulnerable. Second, for this
12 vulnerable reverse-KL case, we study whether methods designed to stabilize asyn-
13 chronous RL can mitigate OPD staleness. We find that they do not improve over
14 a simpler OPD-specific surrogate: recomputing the reverse-KL signal under the
15 current student at learner time without clipping. Third, we identify an OPD-specific
16 teacher-cache constraint: under asynchronous execution, teacher scores are avail-
17 able at learner time only on cached actions. The resulting bias-variance tradeoff
18 for sparse and sampled reverse-KL OPD implementations motivates multi-sample
19 Monte Carlo (MC), which preserves MC correctness while reducing one-sample
20 variance. Finally, we present and open-source **AsyncOPD**, a fully asynchronous
21 OPD training pipeline built from these estimator choices. Experiments show that
22 AsyncOPD improves training throughput by $1.6\times$ to $3.8\times$ over strict synchronous
23 training while reaching comparable accuracy.¹

24 1 Introduction

25 On-policy distillation (OPD) [1, 6, 20] and reinforcement learning (RL) [26, 28] have become
26 central post-training methods for improving large language model (LLM) reasoning [7], including
27 mathematics [17] and coding [27]. OPD trains a student on its own rollouts using dense token-level
28 feedback from a teacher [12], whereas RL learns from reward feedback on rollouts. OPD provides
29 an effective and efficient route for LLM post-training, especially for smaller student models [24].
30 Recent work shows that OPD is not limited to distilling large teachers into small students: it
31 also supports on-policy self-distillation [32] and multi-teacher distillation from domain-specialized
32 teachers comparable in size to the student [2, 21].

33 OPD and RL inherit an on-policy systems bottleneck: each learner update must wait for fresh rollouts
34 from the model being trained [5]. For reasoning tasks, these rollouts are long and expensive, so
35 synchronous training often waits on generation rather than updating the model, leaving learners
36 underutilized. Asynchronous RL [4, 14] relieves this bottleneck by decoupling rollout generation
37 from learner updates: rollout workers keep generating data while the learner updates on earlier

¹Code is available at <https://anonymous.4open.science/r/async-opd-1D8D>.

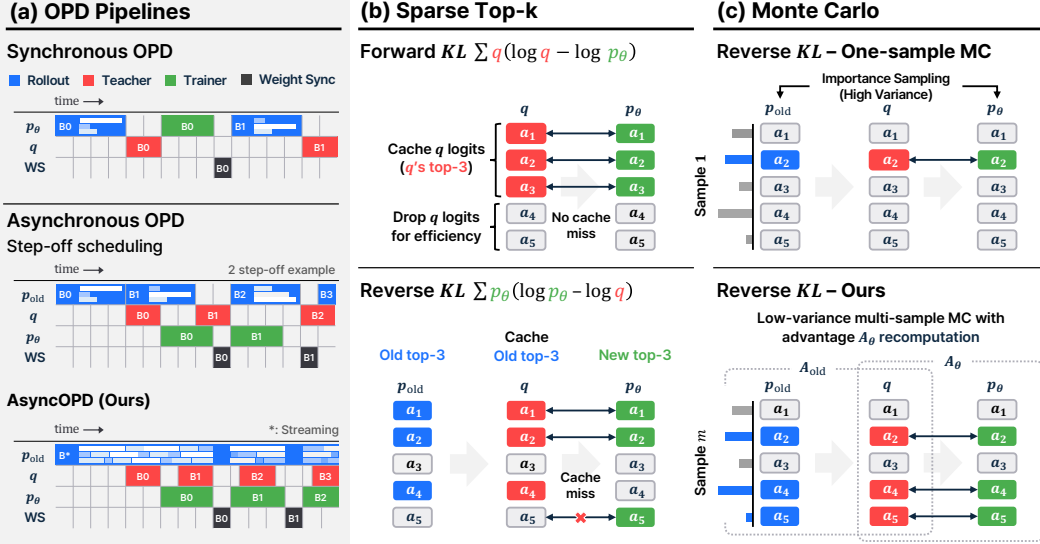


Figure 1: AsyncOPD and stale-cache estimator design. (a) Synchronous OPD is barriered; step-off scheduling [19] overlaps stages but keeps gated rollout batches, while our AsyncOPD streams rollout to reduce long-tail waiting. (b) Sparse top- k exposes a support mismatch under staleness: forward KL is teacher-supported, but reverse KL is student-supported and may require actions outside the cached teacher-scored support. (c) One-sample Monte Carlo is correctable in expectation by importance sampling, but has high variance; our estimator recomputes A_{θ} at learner time and uses multi-sample MC to reduce variance.

38 rollouts, improving training efficiency and hardware utilization [18, 23, 34]. A similar pipeline can
 39 be applied to OPD by running student rollout, teacher scoring, and learner updates in parallel [19].

40 However, asynchronous execution introduces stale-policy data, and learning from such data can
 41 degrade model quality [3]. This creates a trade-off: more aggressive asynchrony improves training
 42 throughput, but it also increases the policy lag between rollout and learning. Prior work on
 43 asynchronous RL therefore studies how to stabilize learning from stale-policy data [4, 10, 33], but
 44 OPD adds a constraint that RL does not have: teacher scores are expensive and are usually cached
 45 only on a finite set of actions. Once the learner receives a teacher-scored cache, it can recompute
 46 current-student log probabilities on cached actions, but it cannot recover teacher scores for actions
 47 that were never scored. This raises three questions that structure our study: (i) how asynchronous
 48 OPD behaves under staleness, (ii) whether asynchronous RL ideas and stale-data solutions transfer to
 49 OPD, and (iii) how OPD-specific constraints shape estimator design.

50 First, we study how KL direction shapes staleness. Under the asynchronous OPD setting, the same
 51 stale rollout cache can affect different KL objectives differently. As illustrated in Fig. 1, forward KL
 52 is teacher-weighted and remains robust to stale rollouts, whereas reverse KL is student-weighted and
 53 becomes vulnerable when current-student actions fall outside the scored cache. We therefore focus
 54 on reverse-KL OPD in the remainder of the staleness analysis.

55 Second, focusing on the reverse-KL case, we ask whether methods designed to stabilize asynchronous
 56 RL can also mitigate OPD staleness. This comparison is natural because reverse KL in OPD admits an
 57 RL-style policy-gradient surrogate, where the teacher-student log-ratio acts as a token-level advantage.
 58 We therefore evaluate PPO-style clipping [16], decoupled PPO [4], and M2PO [33]. We find that
 59 they do not improve over a simpler OPD-specific surrogate: recomputing the reverse-KL token-level
 60 advantage under the current student at learner time without clipping.

61 Third, we return to the teacher-cache constraint and study the resulting bias-variance tradeoff for
 62 sparse and sampled reverse-KL OPD implementations. Stale student top- k supports provide deter-
 63 ministic coverage but are support-mismatched because they may omit actions required by the current
 64 top- k objective, and reweighting inside the stale support cannot recover the missing teacher scores.
 65 One-sample Monte Carlo (MC) avoids this fixed-support mismatch through importance-correctable

66 samples from the stale rollout policy, but suffers from high variance. This motivates multi-sample
67 MC, which caches and teacher-scores multiple stale-policy samples at each decoding step, preserving
68 MC correctability while reducing one-sample variance.

69 Finally, we instantiate these findings in **AsyncOPD**, a fully asynchronous OPD pipeline that overlaps
70 student rollout, teacher scoring, and learner updates. On Qwen3-Base models, AsyncOPD improves
71 training throughput by $1.6\times$ to $3.8\times$ over strict synchronous training while maintaining comparable
72 accuracy. Our contributions are:

- 73 • We provide the first systematic study of staleness in asynchronous OPD through the lens of an
74 OPD-specific teacher-cache constraint.
- 75 • We show that KL direction changes the stale-data problem: forward KL is comparatively robust to
76 stale rollouts, whereas reverse KL is vulnerable because it is student-weighted (Section 4).
- 77 • We identify that the most effective reverse-KL policy-gradient surrogate uses the advantage re-
78 computed at learner time without clipping, and that advanced asynchronous RL surrogates do not
79 improve over this choice (Section 5).
- 80 • We show that stale student top- k supports are support-mismatched, while one-sample MC remains
81 correctable but high-variance; this motivates multi-sample MC (Section 6).
- 82 • We present and open-source **AsyncOPD**, a fully asynchronous OPD training pipeline, and demon-
83 strate improved training efficiency while maintaining OPD quality (Section 7).

84 2 Related Works

85 **On-Policy Distillation** On-policy distillation (OPD) trains a student on its own rollouts while
86 using a teacher to provide dense token-level feedback on the visited prefixes [12, 20]. GKD [1]
87 introduced a token-level KL formulation, while MiniLLM [6] studied a sequence-level reverse-KL
88 variant. Li et al. [11] study token-level OPD training dynamics and recipes for unstable configurations.
89 TIP [22] characterizes per-token importance through student entropy and teacher-student divergence.
90 G-OPD [25] interprets token-level OPD as dense KL-constrained RL and extends it with reward
91 scaling. These works clarify OPD as an effective post-training objective, but assume rollouts, teacher
92 scoring, and learner updates stay synchronized.

93 **Asynchronous RL** In synchronous RL pipelines, training often waits for the longest rollout in
94 a batch to finish, leaving learner resources idle. Asynchronous RL improves hardware utilization
95 by decoupling rollout generation from learner updates. Async RLHF [14] overlaps generation and
96 learning so that new samples are produced while the learner trains on earlier ones. StreamRL [34]
97 further disaggregates the RLHF pipeline into streaming stages. AReaL [4] fully decouples rollout
98 workers from training workers for continuous asynchronous execution. Laminar [18] uses fine-grained
99 weight synchronization for trajectory-level asynchrony. However, asynchronous RL must learn from
100 stale-policy data. Decoupled PPO [4] stabilizes asynchronous RL training by separating the behavior
101 policy for stale rollouts from the proximal policy that anchors PPO [16] updates. M2PO [33] stabilizes
102 stale updates with second-moment importance-weight constraints, and A-3PO [10] reduces decoupled
103 PPO overhead through staleness-aware interpolation.

104 **Asynchronous OPD** VeRL [19] implements step-off OPD schedulers that overlap student rollout,
105 teacher scoring, and learner update by fixing rollout lag to one or two learner steps. These schedulers
106 establish the practical feasibility of asynchronous OPD, but leave open how OPD estimators behave
107 under stale teacher-scored caches. KDFlow [29] improves systems efficiency for LLM distillation
108 by decoupling teacher inference from learner training and transmitting teacher hidden states, but
109 targets synchronous OPD and leaves asynchronous execution as future work. We study this missing
110 asynchronous OPD regime directly and build AsyncOPD from the resulting estimator choices.

111 3 Preliminaries: On-Policy Distillation

112 **OPD setup** At each decoding timestep, we view the visited prefix s as the local state and the next
113 token a as the action. Let $q(a | s)$ denote the teacher policy and $p_\theta(a | s)$ denote the current student
114 policy. Following prior work on token-level OPD [11, 12, 25], we apply local losses to generated

115 output tokens and analyze the resulting objectives at a fixed prefix state s . OPD can be defined with
 116 different divergences; forward and reverse KL are two standard choices [1].

117 **Forward-KL OPD** At a fixed prefix s , forward-KL OPD is teacher-weighted:

$$D_F(\theta; s) = \text{KL}(q(\cdot | s) \| p_\theta(\cdot | s)) = \sum_{a \in \mathcal{V}} q(a | s) (\log q(a | s) - \log p_\theta(a | s)). \quad (1)$$

118 The local gradient is supervised-learning-like: $\nabla_\theta D_F(\theta; s) = - \sum_{a \in \mathcal{V}} q(a | s) \nabla_\theta \log p_\theta(a | s)$.

119 **Reverse-KL OPD** At the same prefix, reverse-KL OPD is student-weighted:

$$D_R(\theta; s) = \text{KL}(p_\theta(\cdot | s) \| q(\cdot | s)) = - \sum_{a \in \mathcal{V}} p_\theta(a | s) (\log q(a | s) - \log p_\theta(a | s)). \quad (2)$$

120 Differentiating and using $\mathbb{E}_{a \sim p_\theta(\cdot | s)}[\nabla_\theta \log p_\theta(a | s)] = \nabla_\theta \sum_a p_\theta(a | s) = 0$ gives

$$\begin{aligned} \nabla_\theta D_R(\theta; s) &= \sum_{a \in \mathcal{V}} p_\theta(a | s) (\log p_\theta(a | s) - \log q(a | s) + 1) \nabla_\theta \log p_\theta(a | s) \\ &= - \sum_{a \in \mathcal{V}} p_\theta(a | s) (\log q(a | s) - \log p_\theta(a | s)) \nabla_\theta \log p_\theta(a | s). \end{aligned} \quad (3)$$

121 Viewing Eq. (3) as a policy-gradient estimator and $A = \log q(a | s) - \log p(a | s)$ as the advantage
 122 term connects reverse-KL OPD to standard RL training machinery, and practical implementations
 123 typically use PPO-style surrogates. Given behavior-policy samples $a \sim p_{\text{beh}}$, define $\rho_\theta(a, s) =$
 124 $p_\theta(a | s) / p_{\text{beh}}(a | s)$ and $\bar{\rho}_\theta(a, s) = \text{clip}(\rho_\theta(a, s), 1 - \epsilon, 1 + \epsilon)$. The PPO-style local surrogate uses
 125 these ratios with a frozen behavior-time signal $A_{\text{beh}}(a, s)$, where $\text{sg}(\cdot)$ denotes stop-gradient:

$$L_{\text{PPO}}(\theta; A_{\text{beh}}) = - \mathbb{E}_{a \sim p_{\text{beh}}} [\min(\rho_\theta \text{sg}(A_{\text{beh}}), \bar{\rho}_\theta \text{sg}(A_{\text{beh}}))]. \quad (4)$$

126 **Sparse and sampled implementations** The dense objectives above are full-vocabulary references.
 127 Practical OPD instead evaluates local KL losses on finite supports or sampled actions [11, 12], trading
 128 computation against support coverage and estimator variance. **Sparse top- k** implementations choose
 129 a support $S(s)$ and evaluate the corresponding restricted KL after renormalizing teacher and student
 130 distributions on $S(s)$. **Monte Carlo (MC)** implementations draw actions from a proposal distribution
 131 and estimate the corresponding local gradient; for reverse KL, this yields the student-sampled
 132 policy-gradient estimator. Details are in Section A.

133 4 Forward- and Reverse-KL OPD Under Staleness

134 Asynchronous OPD has both prefix-level and action-level staleness. Once a rollout is generated, its
 135 visited prefixes are fixed, so an action-level estimator cannot change which states the learner sees.
 136 We therefore focus on the action-level staleness that estimator design can directly address.

137 4.1 Asynchronous OPD Setup

138 Asynchronous OPD is a cached-data pipeline: rollout first selects prefixes and actions, teacher scoring
 139 then annotates those actions, and the learner updates the student later. Unlike synchronous OPD,
 140 these stages are separated in time, so the visited prefixes, action cache, teacher scores, and update
 141 policy may be tied to different student versions. Figure 1 summarizes this cached-teacher setting and
 142 the estimator contrasts induced by the three-stage cache pipeline.

143 **Teacher-cache constraint** Full-vocabulary teacher logits allow dense KL computation, but caching
 144 and transferring them is prohibitively expensive, especially in an asynchronous pipeline. We therefore
 145 focus on sparse top- k supports and MC samples as the sparse and sampled cases.

146 **Stage 1: Student rollout** A rollout actor samples trajectories from a stale student p_{old} , which fixes
 147 the visited prefixes s . At each prefix, it stores cached actions $C_{\text{old}}(s)$ together with their rollout-time
 148 log probabilities under p_{old} , such as a sampled token or a top- k support.

149 **Stage 2: Teacher scoring** Let $C_{\text{score}}(s)$ denote the teacher-scored cache; it may come from the
 150 rollout cache $C_{\text{old}}(s)$ or be selected by the teacher at scoring time. Once teacher scoring is complete,
 151 teacher logits are available only on $C_{\text{score}}(s)$.

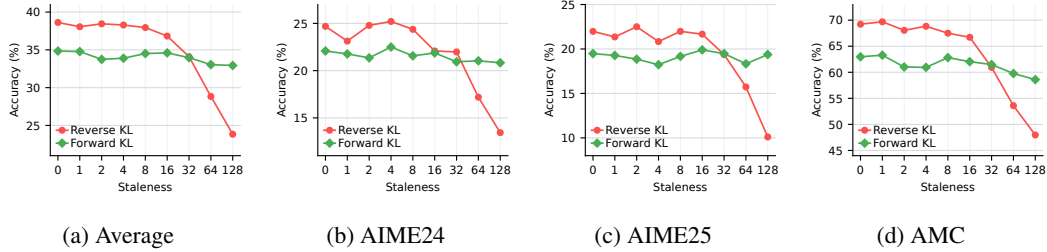


Figure 2: Accuracy comparison under staleness for forward- and reverse-KL OPD. Reverse KL starts higher at zero staleness but degrades faster as staleness grows; forward KL is flatter across the sweep.

152 **Stage 3: Student update** By learner update time, the student has moved to the current policy p_θ .
 153 The learner can recompute $\log p_\theta(a | s)$ for $a \in C_{\text{score}}(s)$, but the current local OPD objective may
 154 place mass on actions outside this teacher-scored cache, such as the current student top- k support or
 155 current student-sampled actions. Thus the learner can update current student probabilities on cached
 156 actions, but cannot recover teacher signals for missing actions without additional teacher access.

157 **Experimental setup** Unless otherwise stated, we train a Qwen3-4B-Base student using a Qwen3-
 158 30B-A3B-Instruct-2507 teacher [24]. The training data is DeepMath [8], filtered to 57,630 math
 159 problems with difficulty level at least 6, and we report final-checkpoint Avg@32 accuracy on
 160 AIME24 [30], AIME25 [31], and AMC [13]. Experimental details are provided in Section B; dataset
 161 and metric details are in Section C.

162 4.2 Forward KL vs. Reverse KL Under Staleness

163 The KL direction fixes the action weighting: forward KL weights actions by the teacher q , whereas
 164 reverse KL weights them by the student p_θ . With cached teacher scores, this weighting difference
 165 becomes a support-ownership difference (Fig. 1). Under a scored-cache restriction, this makes
 166 forward KL less exposed to stale student action choices: it does not need to convert stale student-
 167 sampled actions into a current-student expectation. Reverse KL instead depends on student-weighted
 168 action terms, so the same asynchronous cache creates a different action-level staleness problem.

169 **Experimental results** Figure 2 compares representative practical OPD implementations from prior
 170 work: sparse top- k forward KL [19] and PPO-style reverse-KL surrogates [12, 25]. Reverse KL
 171 starts higher at zero staleness, but as staleness increases it drops faster and is eventually overtaken by
 172 forward KL. We therefore focus the rest of the staleness analysis on how to make reverse-KL OPD
 173 robust under larger rollout staleness.

174 **Finding 1.** Forward KL is teacher-weighted and robust to rollout staleness, whereas reverse KL
 is student-weighted and vulnerable to rollout staleness.

175 **Two axes of reverse-KL staleness** The cache analysis above suggests a possible mechanism
 176 for this gap: because reverse KL is weighted by the current student, stale teacher-scored caches
 177 may fail to cover actions needed by the current reverse-KL objective. In addition, reverse-KL
 178 policy-gradient updates can be instantiated with multiple stale-data surrogates, including PPO-style
 179 and asynchronous-RL variants. We therefore split the reverse-KL analysis into a **policy-gradient**
 180 **surrogate axis**, studied in Section 5, and a **cached-support axis**, studied in Section 6.

181 5 Reverse-KL: Policy-Gradient Surrogates Under Staleness

182 Reverse-KL OPD admits several policy-gradient surrogate choices under stale rollouts. This section
 183 compares which choices remain effective under staleness.

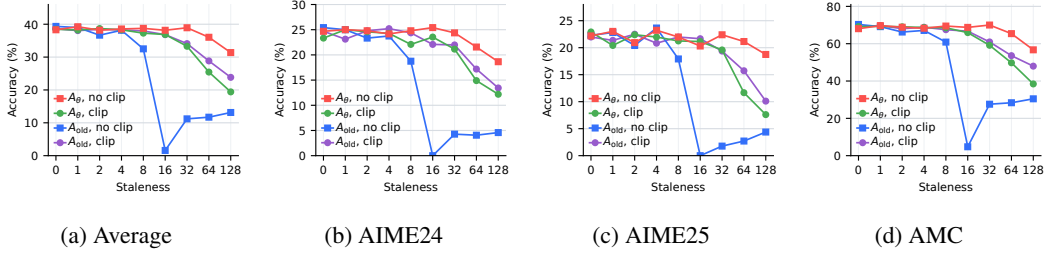


Figure 3: Accuracy comparison under staleness for the advantage-and-clipping ablation. Recomputing A_θ at learner time and avoiding clipping gives the most stable performance across the sweep, while clipping mainly helps the frozen A_{old} baseline.

184 5.1 Policy-Gradient Surrogate Choices

185 **PPO-style objective** In the PPO-style surrogate in Eq. (4), the advantage is computed under the
 186 behavior policy and then held fixed during the learner update. In stale reverse-KL OPD, the behavior
 187 policy is the rollout student, so a mechanical PPO-style adaptation sets $p_{\text{beh}} = p_{\text{old}}$ and uses the
 188 rollout-time reverse-KL advantage $A_{\text{old}}(a, s) = \log q(a | s) - \log p_{\text{old}}(a | s)$ as A_{beh} , together with
 189 the clipped old-to-current ratio. The unclipped variant simply drops the clipped term.

190 **Exact importance-sampling identity** In contrast, rewriting the reverse-KL objective (Eq. (2)) by
 191 importance sampling suggests a different surrogate choice. With the current reverse-KL advantage
 192 $A_\theta(a, s) = \log q(a | s) - \log p_\theta(a | s)$, and assuming p_{old} has support wherever p_θ does, the current
 193 reverse-KL objective admits the exact old-to-current importance-sampling (IS) identity

$$D_R(\theta; s) = -\mathbb{E}_{a \sim p_\theta} [A_\theta(a, s)] = -\mathbb{E}_{a \sim p_{\text{old}}} [\rho_\theta(a, s) A_\theta(a, s)]. \quad (5)$$

194 For the policy-gradient update, the advantage is used as a stop-gradient weight; the derivative of the
 195 omitted A_θ term cancels by the score-function identity, as in Eq. (3). Thus the IS view points to the
 196 opposite surrogate choice from the mechanical PPO adaptation: recompute A_θ under the current
 197 student and use the old-to-current ratio without clipping, with A_θ treated as a stop-gradient advantage.

198 **A two-by-two surrogate ablation** The PPO-style adaptation and the OPD/IS identity suggest
 199 different surrogate choices. We therefore ablate the advantage (A_{old} versus A_θ) and whether to clip
 200 the ratio, with $\text{sg}(\cdot)$ denoting stop-gradient:

$$L_{\text{old}}^{\text{clip}}(\theta) = -\mathbb{E}_{a \sim p_{\text{old}}} [\min(\rho_\theta \text{sg}(A_{\text{old}}), \bar{\rho}_\theta \text{sg}(A_{\text{old}}))], \quad L_{\text{old}}^{\text{noclip}}(\theta) = -\mathbb{E}_{a \sim p_{\text{old}}} [\rho_\theta \text{sg}(A_{\text{old}})], \quad (6)$$

$$L_\theta^{\text{clip}}(\theta) = -\mathbb{E}_{a \sim p_{\text{old}}} [\min(\rho_\theta \text{sg}(A_\theta), \bar{\rho}_\theta \text{sg}(A_\theta))], \quad L_\theta^{\text{noclip}}(\theta) = -\mathbb{E}_{a \sim p_{\text{old}}} [\rho_\theta \text{sg}(A_\theta)]. \quad (7)$$

201 Here $L_{\text{old}}^{\text{clip}}$ is the PPO-style adaptation, while L_θ^{noclip} is the OPD/IS surrogate.

202 **Advanced asynchronous RL surrogates** Decoupled PPO [4] and M2PO [33] are asynchronous
 203 RL surrogates designed to improve robustness to stale-policy updates. We evaluate whether these
 204 previously unstudied asynchronous RL surrogates also help OPD under staleness.

205 5.2 Experimental Results

206 Fig. 3 and Table 1a compare the four combinations of A_{old} versus A_θ and
 207 clipping versus no clipping. The best variant is the OPD/IS choice: A_θ
 208 without clipping. The PPO-style baseline, A_{old} with clipping, remains
 209 a strong stale-surrogate baseline. Clipping helps A_{old} by limiting stale,
 210 large-ratio updates, but hurts A_θ : recomputing A_θ already reduces the
 211 high-percentile ρ_θ tail at staleness 64 (Fig. 4), so clipping removes useful
 212 signal. Likewise, Fig. 5 and Table 1a show that advanced asynchronous
 213 RL surrogates such as decoupled PPO and M2PO do not outperform A_θ
 214 without clipping, which becomes our reference surrogate below.

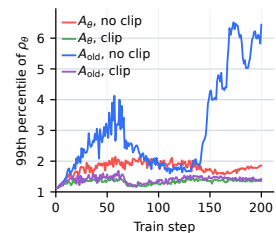


Figure 4: A_θ reduces the p99 ρ_θ tail under no clip.

215 **Finding 2.** The most effective reverse-KL correction is to recompute A_θ at learner time without clipping; advanced asynchronous RL surrogates such as decoupled PPO and M2PO do not improve over it.

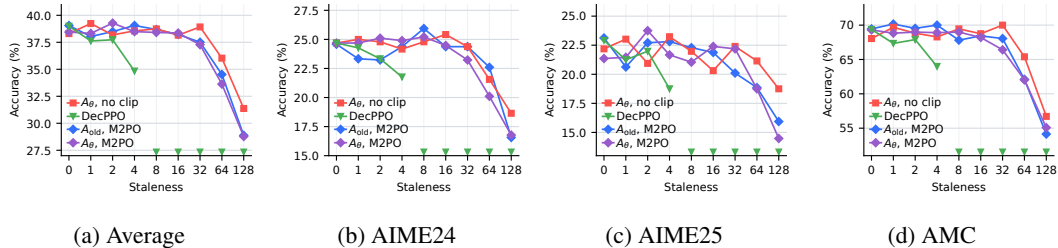


Figure 5: Accuracy comparison under staleness for advanced asynchronous RL surrogates. Decoupled PPO [4] and M2PO [33] do not consistently improve over the simpler OPD/IS surrogate that recomputes A_θ without clipping; Decoupled PPO is clipped for readability because of low accuracy.

216 6 Reverse-KL: Cached Supports Under Staleness

217 Having fixed A_θ without clipping as the reference surrogate, we now ask which cached actions
 218 provide the teacher scores needed to evaluate it, and how to improve this cached-support estimator.
 219 This cached-support axis is specific to OPD because teacher scoring is local and expensive: the
 220 teacher cache determines which actions have teacher scores available to the learner.

221 **Sparse top- k : stale-support biased** Although sparse top- k is biased relative to the dense reverse-
 222 KL objective, it is a practical low-variance approximation on the current student support $S_\theta(s) =$
 223 $\text{TopK}(p_\theta(\cdot | s), k)$. Under asynchronous rollout reuse, however, teacher scores are cached on the
 224 rollout-time support $S_{\text{old}}(s) = \text{TopK}(p_{\text{old}}(\cdot | s), k)$, which may miss actions in the current support
 225 $S_\theta(s)$. Reweighting within S_{old} cannot recover these missing teacher scores, so stale sparse top- k
 226 remains a support-mismatched approximation, not an exact correction of the current top- k objective.

227 **One-sample MC: correctable but high variance** Sampled-token MC instead caches an action
 228 drawn from a behavior distribution: $a \sim p_{\text{old}}(\cdot | s)$ together with $\log p_{\text{old}}(a | s)$. When the behavior
 229 policy covers the current policy support, exact old-to-current IS gives an unbiased estimator of
 230 the current reverse-KL fixed-prefix gradient. Thus one-sample MC is action-level correctable in
 231 expectation, but the resulting IS estimator can have high variance. This proposal-sampling structure
 232 is the key contrast with stale top- k , whose actions come from a deterministic stale support.

233 6.1 Proposed Solution: Multi-Sample MC

234 **Multi-sample MC: correctable with reduced variance** We propose multi-sample MC to reduce
 235 the high variance of one-sample MC by caching multiple behavior-policy samples at the same prefix
 236 and averaging their IS-corrected local gradients.

237 Multi-sample MC is especially natural in asynchronous OPD. In RL for LLM post-training, branching
 238 a prefix into multiple actions is expensive because each branch typically requires a full continuation
 239 before the reward or advantage can be evaluated. In synchronous OPD, sparse top- k already provides
 240 a low-variance approximation and one-sample MC provides an unbiased sampled gradient estimator,
 241 so there is little motivation to cache multiple sampled actions per prefix. Under asynchronous OPD,
 242 this tradeoff changes: sparse top- k becomes the stale fixed-support approximation analyzed above,
 243 and one-sample MC remains correctable but high-variance, making multi-sample MC a natural
 244 cached-support estimator for asynchronous OPD.

245 Concretely, for each visited prefix s , rollout samples $a_1, \dots, a_m \sim p_{\text{old}}(\cdot | s)$ and caches their rollout
 246 log probabilities and teacher scores. At learner time, we recompute $A_\theta(a_i, s)$ and use the averaged
 247 unclipped old-to-current IS surrogate $\hat{L}_m^{\text{MC}}(\theta; s) = -\frac{1}{m} \sum_{i=1}^m \rho_\theta(a_i, s) \text{sg}(A_\theta(a_i, s))$. By linearity,
 248 the gradient has the same expectation as the one-sample MC estimator; averaging independent
 249 behavior-policy samples reduces the Monte Carlo variance.

Table 1: Staleness-sensitivity slopes. Entries fit accuracy against $\log_2(\text{staleness} + 1)$; more negative values indicate stronger degradation with staleness.

(a) Policy-gradient surrogates							(b) Multi-sample MC				
Benchmark	A_{old} (clip)	A_{old}	A_{θ} (clip)	A_{θ}	A_{old} (M2PO)	A_{θ} (M2PO)	Benchmark	MC1	MC4	MC16	MC64
AIME24	-1.44	-3.99	-1.64	-0.69	-0.66	-1.00	AIME24	-0.69	<u>-0.42</u>	-0.46	-0.34
AIME25	-1.42	-3.69	-1.88	-0.38	<u>-0.78</u>	-0.81	AIME25	-0.38	<u>-0.23</u>	-0.24	-0.11
AMC	-3.06	-8.00	-4.06	-1.12	-1.75	<u>-1.72</u>	AMC	-1.12	-0.87	-0.74	<u>-0.84</u>

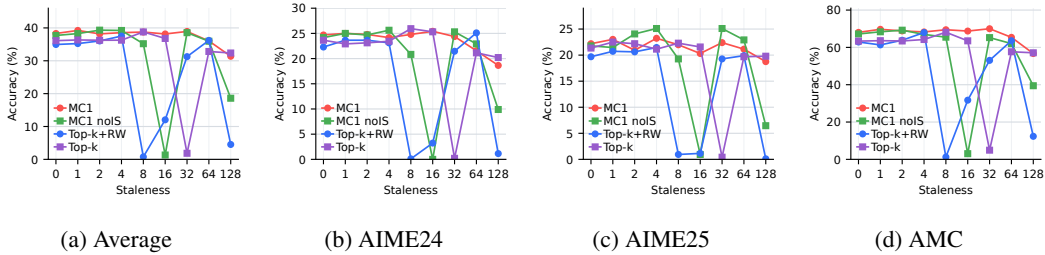


Figure 6: Accuracy comparison under staleness for sampled MC versus stale top- k . Top- k +RW denotes reweighting on the stale top- k support. Old-to-current IS corrects sampled MC in expectation, whereas reweighting cannot repair the missing teacher scores induced by stale top- k supports.

250 6.2 Experimental Results

251 **Sparse top- k vs. one-sample MC** Figure 6 compares one-sample MC and sparse top- k , with and
 252 without old-to-current reweighting. For one-sample MC, IS substantially improves robustness as
 253 staleness increases. For sparse top- k , the same reweighting does not improve performance, since it
 254 cannot recover missing current-support actions. As a result, one-sample MC with IS is the strongest of
 255 the four methods, consistent with the support-correctability analysis above. We include an additional
 256 ablation disentangling MC sample count from IS in Section E.

257 **One-sample MC vs. multi-sample MC** Figure 7 and Table 1b show that multi-sample MC
 258 improves one-sample MC at large staleness: $m = 4$ already gives a clear jump, while $m \in \{4, 16, 64\}$
 259 performs similarly.

260 **Finding 3.** One-sample MC is more effective than stale sparse top- k ; multi-sample MC further
 261 improves this estimator by reducing one-sample variance while preserving MC correctability.

261 7 AsyncOPD: Fully Asynchronous OPD

262 AsyncOPD is our fully asynchronous OPD system. Following AReL [4], it overlaps rollout, teacher
 263 scoring, and learner updates.

264 **Scheduler** The step-off scheduler family was originally implemented in VeRL [19]: a k -step-off
 265 run fixes rollout lag to k learner updates, but still waits for complete rollout batches. AsyncOPD
 266 streams examples instead: workers pause only for weight sync, preserve in-flight prefixes, teacher
 267 scoring consumes completed items, and the learner updates once a scored batch is ready (Fig. 1).

268 **Experimental setup** The main comparison uses Qwen3- $\{1.7B, 4B, 8B\}$ -Base students with the
 269 Qwen3-30B-A3B-Instruct-2507 teacher. All runs use the same reverse-KL estimator: current-policy
 270 A_{θ} , no clipping, old-to-current IS, and either MC64 or MC1. We compare strict sync, two-step-off,
 271 and AsyncOPD for 100 training iterations on the same 8-GPU node; all AsyncOPD runs use $\tau = 4$.
 272 Section F gives GPU allocation, queue-depth, and scheduler details.

273 **Experimental Results** Table 2 reports training throughput, pipeline overlap (average concurrent
 274 OPD-stage activity), and final AIME24 Avg@32 for the Qwen3-Base students. AsyncOPD achieves

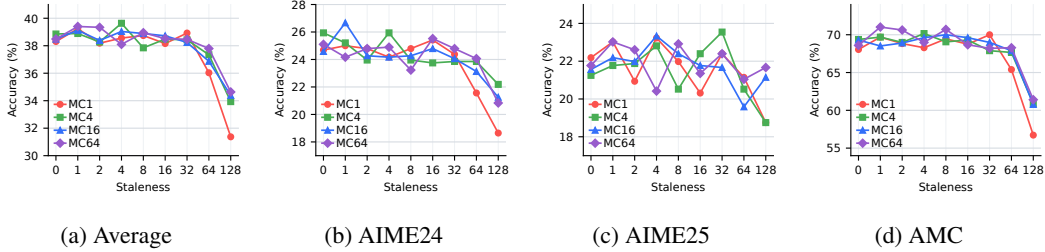


Figure 7: Accuracy comparison under staleness for multi-sample MC. Increasing the number of samples improves large-staleness behavior.

Table 2: AsyncOPD scheduler results for Qwen3-Base models. Train tok/s is training throughput; parentheses show speedup over the matched strict-sync baseline. Overlap is concurrent OPD-stage activity. Avg@32 is final AIME24. AsyncOPD achieves the highest throughput and overlap in all matched settings while maintaining comparable final accuracy.

Student	Scheduler	MC64			MC1		
		Train tok/s (\times sync)	Overlap	Avg@32	Train tok/s (\times sync)	Overlap	Avg@32
Qwen3-1.7B-Base	Strict sync	8.7k (1.00 \times)	0.81	8.85	8.6k (1.00 \times)	0.82	8.65
	Two-step-off	14.2k (1.64 \times)	1.49	8.23	18.5k (2.15 \times)	1.65	8.12
	AsyncOPD (ours)	23.4k (2.70\times)	2.13	9.38	28.1k (3.28\times)	2.31	8.44
Qwen3-4B-Base	Strict sync	9.5k (1.00 \times)	0.81	25.00	8.1k (1.00 \times)	0.82	23.33
	Two-step-off	12.4k (1.30 \times)	1.64	23.85	12.9k (1.60 \times)	1.65	23.54
	AsyncOPD (ours)	15.8k (1.66\times)	2.19	25.00	16.4k (2.03\times)	2.27	23.44
Qwen3-8B-Base	Strict sync	7.5k (1.00 \times)	0.81	26.56	6.5k (1.00 \times)	0.84	28.44
	Two-step-off	11.6k (1.55 \times)	1.76	26.56	9.6k (1.47 \times)	1.63	28.85
	AsyncOPD (ours)	14.5k (1.94\times)	2.24	28.65	10.6k (1.63\times)	2.17	27.50

275 the highest throughput and overlap in every matched comparison. In MC64, it reaches up to 2.7 \times the
 276 strict-sync throughput while achieving the best or tied-best final accuracy. MC1 shows the same trend:
 277 AsyncOPD delivers the highest throughput (up to 3.3 \times strict-sync) and overlap for every student,
 278 with competitive final accuracy. Train-time accuracy curves are reported in Section F.

279 8 Conclusion

280 We present the first systematic study of staleness in asynchronous on-policy distillation (OPD). Our
 281 results show that KL direction shapes the stale-data problem: forward KL remains robust to stale
 282 rollouts, whereas reverse KL is more vulnerable because it is student-weighted. In reverse-KL OPD,
 283 the most effective policy-gradient surrogate uses the current advantage recomputed at learner time
 284 without clipping; advanced asynchronous RL surrogates do not improve over this choice. We also
 285 find that stale student top- k supports are support-mismatched, whereas one-sample Monte Carlo (MC)
 286 remains correctable but high-variance. This contrast motivates multi-sample MC, which preserves MC
 287 correctability while reducing one-sample variance. Finally, we present and open-source **AsyncOPD**,
 288 a fully asynchronous OPD training pipeline built from these estimator choices, improving training
 289 efficiency while maintaining OPD quality.

290 **Limitations and Future Work** We study sparse and Monte Carlo OPD estimators, not dense
 291 full-vocabulary KL in the asynchronous setting. Although dense KL avoids cached-support mismatch,
 292 it is difficult to implement efficiently when rollout, teacher scoring, and learner updates are decoupled.
 293 KDFlow [29] suggests one path by transmitting teacher hidden states and recomputing student logits,
 294 but only for synchronous OPD. Extending this approach to asynchronous OPD while handling stale
 295 rollouts and preserving throughput is an important future direction. Our experiments are also limited
 296 to a single 8-GPU node by available resources, not by the pipeline itself; scaling to larger multi-node
 297 clusters remains future work.

298 **References**

- 299 [1] R. Agarwal, N. Vieillard, Y. Zhou, P. Stanczyk, S. R. Garea, M. Geist, and O. Bachem. On-
300 policy distillation of language models: Learning from self-generated mistakes. In *The Twelfth*
301 *International Conference on Learning Representations*, 2024. URL [https://openreview.](https://openreview.net/forum?id=3zKtaqxLhW)
302 [net/forum?id=3zKtaqxLhW](https://openreview.net/forum?id=3zKtaqxLhW).
- 303 [2] DeepSeek-AI. Deepseek-v4: Towards highly efficient million-token context intelligence, 2026.
- 304 [3] F. Devvrit, L. Madaan, R. Tiwari, R. Bansal, S. S. Duvvuri, M. Zaheer, I. S. Dhillon,
305 D. Brandfonbrener, and R. Agarwal. The art of scaling reinforcement learning compute
306 for LLMs. In *The Fourteenth International Conference on Learning Representations*, 2026.
307 URL <https://openreview.net/forum?id=FMjeC9Msws>.
- 308 [4] W. Fu, J. Gao, X. Shen, C. Zhu, Z. Mei, C. He, S. Xu, G. Wei, J. Mei, W. JIASHU, T. Yang,
309 B. Yuan, and Y. Wu. AREAL: A large-scale asynchronous reinforcement learning system for
310 language reasoning. In *The Thirty-ninth Annual Conference on Neural Information Processing*
311 *Systems*, 2025. URL <https://openreview.net/forum?id=X9diEuva9R>.
- 312 [5] W. Gao, Y. Zhao, D. An, T. Wu, L. Cao, S. Xiong, J. Huang, W. Wang, S. Yang, W. Su, et al.
313 Rollpacker: Mitigating long-tail rollouts for fast, synchronous rl post-training. *arXiv preprint*
314 *arXiv:2509.21009*, 2025.
- 315 [6] Y. Gu, L. Dong, F. Wei, and M. Huang. MiniLLM: Knowledge distillation of large language
316 models. In *The Twelfth International Conference on Learning Representations*, 2024. URL
317 <https://openreview.net/forum?id=5h0qf7IBZZ>.
- 318 [7] D. Guo, D. Yang, H. Zhang, J. Song, P. Wang, Q. Zhu, R. Xu, R. Zhang, S. Ma, X. Bi, et al.
319 Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv*
320 *preprint arXiv:2501.12948*, 2025.
- 321 [8] Z. He, T. Liang, J. Xu, Q. Liu, X. Chen, Y. Wang, L. Song, D. Yu, Z. Liang, W. Wang, et al.
322 Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical
323 dataset for advancing reasoning. *arXiv preprint arXiv:2504.11456*, 2025.
- 324 [9] W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. E. Gonzalez, H. Zhang, and
325 I. Stoica. Efficient memory management for large language model serving with PagedAttention.
326 In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- 327 [10] X. Li, S. Wu, and Z. Shen. A-3po: Accelerating asynchronous llm training with staleness-aware
328 proximal policy approximation. *arXiv preprint arXiv:2512.06547*, 2025.
- 329 [11] Y. Li, Y. Zuo, B. He, J. Zhang, C. Xiao, C. Qian, T. Yu, H.-a. Gao, W. Yang, Z. Liu, et al.
330 Rethinking on-policy distillation of large language models: Phenomenology, mechanism, and
331 recipe. *arXiv preprint arXiv:2604.13016*, 2026.
- 332 [12] K. Lu and T. M. Lab. On-policy distillation. *Thinking Machines Lab: Connectionism*, 2025.
333 doi: 10.64434/tml.20251026. <https://thinkingmachines.ai/blog/on-policy-distillation>.
- 334 [13] Mathematical Association of America. American Mathematics Competitions – AMC. [https://](https://maa.org/)
335 maa.org/, 2023. Accessed 2026-04-03.
- 336 [14] M. Noukhovitch, S. Huang, S. Xhonneux, A. Hosseini, R. Agarwal, and A. Courville. Faster,
337 more efficient RLHF through off-policy asynchronous learning. In *The Thirteenth International*
338 *Conference on Learning Representations*, 2025. URL [https://openreview.net/forum?](https://openreview.net/forum?id=FhTAG591Ve)
339 [id=FhTAG591Ve](https://openreview.net/forum?id=FhTAG591Ve).
- 340 [15] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison,
341 L. Antiga, and A. Lerer. Automatic differentiation in PyTorch. In *NIPS-W*, 2017.
- 342 [16] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
343 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

- 344 [17] Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. Li, Y. Wu, et al.
345 Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv*
346 *preprint arXiv:2402.03300*, 2024.
- 347 [18] G. Sheng, Y. Tong, B. Wan, W. Zhang, C. Jia, X. Wu, Y. Wu, X. Li, C. Zhang, Y. Peng, et al.
348 Laminar: A scalable asynchronous rl post-training framework. *arXiv preprint arXiv:2510.12633*,
349 2025.
- 350 [19] G. Sheng, C. Zhang, Z. Ye, X. Wu, W. Zhang, R. Zhang, Y. Peng, H. Lin, and C. Wu. Hybridflow:
351 A flexible and efficient rlhf framework. In *Proceedings of the Twentieth European Conference*
352 *on Computer Systems*, pages 1279–1297, 2025.
- 353 [20] M. Song and M. Zheng. A survey of on-policy distillation for large language models. *arXiv*
354 *preprint arXiv:2604.00626*, 2026.
- 355 [21] B. Xiao, B. Xia, B. Yang, B. Gao, B. Shen, C. Zhang, C. He, C. Lou, F. Luo, G. Wang, et al.
356 Mimo-v2-flash technical report. *arXiv preprint arXiv:2601.02780*, 2026.
- 357 [22] Y. Xu, H. Sang, Z. Zhou, R. He, Z. Wang, and A. Geramifard. Tip: Token importance in
358 on-policy distillation. *arXiv preprint arXiv:2604.14084*, 2026.
- 359 [23] R. Yan, Y. Jiang, T. Wu, J. Gao, Z. Mei, W. Fu, H. Mai, W. Wang, Y. Wu, and B. Yuan.
360 Areal-hex: Accommodating asynchronous rl training over heterogeneous gpus. *arXiv preprint*
361 *arXiv:2511.00796*, 2025.
- 362 [24] A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv, et al.
363 Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- 364 [25] W. Yang, W. Liu, R. Xie, K. Yang, S. Yang, and Y. Lin. Learning beyond teacher: Generalized
365 on-policy distillation with reward extrapolation. *arXiv preprint arXiv:2602.12125*, 2026.
- 366 [26] Q. Yu, Z. Zhang, R. Zhu, Y. Yuan, X. Zuo, Y. Yue, W. Dai, T. Fan, G. Liu, L. Liu, et al. Dapo:
367 An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*,
368 2025.
- 369 [27] A. Zeng, X. Lv, Z. Hou, Z. Du, Q. Zheng, B. Chen, D. Yin, C. Ge, C. Huang, C. Xie, et al.
370 Glm-5: from vibe coding to agentic engineering. *arXiv preprint arXiv:2602.15763*, 2026.
- 371 [28] K. Zhang, Y. Zuo, B. He, Y. Sun, R. Liu, C. Jiang, Y. Fan, K. Tian, G. Jia, P. Li, et al. A survey
372 of reinforcement learning for large reasoning models. *arXiv preprint arXiv:2509.08827*, 2025.
- 373 [29] S. Zhang, X. Zhang, T. Zhang, B. Hu, Y. Chen, and J. Xu. Kdflow: A user-friendly and efficient
374 knowledge distillation framework for large language models. *arXiv preprint arXiv:2603.01875*,
375 2026.
- 376 [30] Y. Zhang and T. Math-AI. AIME 2024. [https://huggingface.co/datasets/](https://huggingface.co/datasets/Maxwell-Jia/AIME_2024)
377 [Maxwell-Jia/AIME_2024](https://huggingface.co/datasets/Maxwell-Jia/AIME_2024), 2024. Hugging Face dataset; accessed 2026-04-03.
- 378 [31] Y. Zhang and T. Math-AI. AIME 2025. [https://huggingface.co/datasets/](https://huggingface.co/datasets/yentinglin/aime_2025)
379 [yentinglin/aime_2025](https://huggingface.co/datasets/yentinglin/aime_2025), 2025. Hugging Face dataset; accessed 2026-04-03.
- 380 [32] S. Zhao, Z. Xie, M. Liu, J. Huang, G. Pang, F. Chen, and A. Grover. Self-distilled reasoner:
381 On-policy self-distillation for large language models. *arXiv preprint arXiv:2601.18734*, 2026.
- 382 [33] H. Zheng, J. Zhao, and B. Chen. Prosperity before collapse: How far can off-policy RL
383 reach with stale data on LLMs? In *The Fourteenth International Conference on Learning*
384 *Representations*, 2026. URL <https://openreview.net/forum?id=IIg15Mwelz>.
- 385 [34] Y. Zhong, Z. Zhang, X. Song, H. Hu, C. Jin, B. Wu, N. Chen, Y. Chen, Y. Zhou, C. Wan, et al.
386 Streamrl: Scalable, heterogeneous, and elastic rl for llms with disaggregated stream generation.
387 *arXiv preprint arXiv:2504.15930*, 2025.

388 A Sparse and Monte Carlo Reverse-KL Implementations

389 A.1 Sparse Top- k Reverse-KL OPD

390 The dense reverse-KL objective in Eq. (2) sums over the full vocabulary. A sparse top- k implementa-
391 tion instead evaluates reverse KL on a finite student support

$$S_\theta(s) = \text{TopK}(p_\theta(\cdot | s), k). \quad (8)$$

392 For any support S , define the restricted normalizers $Z_p^S(s) = \sum_{u \in S} p_\theta(u | s)$ and $Z_q^S(s) =$
393 $\sum_{u \in S} q(u | s)$, and the renormalized distributions

$$\tilde{p}_\theta^S(a | s) = \frac{p_\theta(a | s) \mathbf{1}[a \in S]}{Z_p^S(s)}, \quad \tilde{q}^S(a | s) = \frac{q(a | s) \mathbf{1}[a \in S]}{Z_q^S(s)}. \quad (9)$$

394 The sparse reverse-KL objective is

$$\begin{aligned} D_R^S(\theta; s) &= \text{KL}(\tilde{p}_\theta^S(\cdot | s) \| \tilde{q}^S(\cdot | s)) \\ &= - \sum_{a \in S} \tilde{p}_\theta^S(a | s) (\log \tilde{q}^S(a | s) - \log \tilde{p}_\theta^S(a | s)). \end{aligned} \quad (10)$$

395 In practice, when $S = S_\theta(s)$, we treat the selected top- k support as fixed during the local update.

396 A.2 Monte Carlo Reverse-KL OPD

397 Let $A_\theta(a, s) = \log q(a | s) - \log p_\theta(a | s)$. From Eq. (3), the dense reverse-KL gradient can be
398 written as

$$\nabla_\theta D_R(\theta; s) = -\mathbb{E}_{a \sim p_\theta(\cdot | s)} [A_\theta(a, s) \nabla_\theta \log p_\theta(a | s)]. \quad (11)$$

399 A one-sample current-policy Monte Carlo estimator is therefore

$$\hat{g}_{\text{MC}}(s, a) = -A_\theta(a, s) \nabla_\theta \log p_\theta(a | s), \quad a \sim p_\theta(\cdot | s). \quad (12)$$

400 With m independent samples $a_i \sim p_\theta(\cdot | s)$, the corresponding multi-sample estimator averages the
401 same local term:

$$\hat{g}_m(s) = -\frac{1}{m} \sum_{i=1}^m A_\theta(a_i, s) \nabla_\theta \log p_\theta(a_i | s). \quad (13)$$

402 B Experimental Details

403 This section details the experimental setup. Unless explicitly stated otherwise, experiments use the
404 common setup in Table 3 and report final-checkpoint Avg@32 accuracy. Our implementation uses
405 vLLM [9] for rollout generation and teacher scoring, PyTorch FSDP [15] for learner training, and
406 runs each experiment on a single $8 \times \text{B200}$ node. Individual experiments take roughly 1–12 hours,
407 depending on the setting. Asset URLs, license names, and versions are summarized in Table 5.

408 **Constructing the staleness axis.** The main text uses staleness as an experimental control over
409 how old the cached rollout data is when the learner updates on it. In all staleness plots and tables in
410 Sections 4 to 6, staleness is measured in train-batch steps. One train-batch step is one logical rollout
411 batch consumed by the learner for a training iteration. The sweep value k is therefore the target
412 number of train-batch steps by which the consumed cache trails the current learner; equivalently, it
413 is the target cache depth in logical rollout batches. A value $k = 0$ is synchronous: rollout, teacher
414 scoring, training, and weight synchronization occur in strict sequence. For $k > 0$, the run first
415 generates exactly k rollout batches with the initial student snapshot before the first learner update.
416 Training then consumes the oldest available generated batch; after each learner update and weight
417 synchronization, a new rollout batch is generated with the latest student snapshot whenever needed to
418 restore the target cache depth.

419 This protocol is the operational source of the prefix- and action-level staleness discussed in the main
420 text: the consumed prefixes and cached actions come from an older rollout student, while the update
421 is applied to the current student. For a consumed rollout batch, let t_{roll} be the train-batch index of the

Table 3: Experimental settings.

Setting	Value
Student	Qwen3-4B-Base for main staleness/estimator experiments; AsyncOPD experiments also use Qwen3-{1.7B,4B,8B}-Base and thinking-disabled Qwen3-{1.7B,4B,8B} [24].
Teacher	Qwen3-30B-A3B-Instruct-2507 [24]
Training data	DeepMath [8] filtered to 57,630 math problems with difficulty level at least 6
Prompt / response lengths	2,048 prompt tokens / 16,384 response tokens
Training horizon	200 training iterations; in the common setup each iteration contains four mini-batch optimizer updates
Optimization	batch size 256, mini-batch size 64
Optimizer	AdamW-style optimizer with learning rate 3×10^{-6} , constant schedule, weight decay 0.01
PPO clipping	$\epsilon = 0.2$ for clipped PPO-style surrogates
M2PO budget	0.01
Rollout generation	temperature 1.0, top- $p = 1.0$, no top- k truncation
Evaluation generation	32 samples per problem, temperature 1.0
Implementation	vLLM rollout and teacher scoring; PyTorch FSDP learner training
Hardware	Single $8 \times B200$ node
GPU allocation	One teacher GPU; strict sync time-shares seven GPUs between rollout and training; stale-cache sweeps and AsyncOPD use concurrent rollout and trainer GPU pools

Table 4: Evaluation datasets.

Dataset	Problems
AIME 2024	30
AIME 2025	30
AMC 2023	40

422 student snapshot used for generation and t_{train} be the train-batch index at learner time. The staleness
 423 used in the plots is

$$\Delta_{\text{batch}} = t_{\text{train}} - t_{\text{roll}}.$$

424 All examples in a logical batch share the same rollout snapshot and therefore share the same Δ_{batch} .
 425 Under the controlled cache protocol, Δ_{batch} ramps as $0, 1, 2, \dots$ while the initial cache is drained and
 426 then plateaus at k . Thus a 64-batch target cache depth is plotted as staleness 64. A train-batch step
 427 can contain multiple mini-batch optimizer updates; in the common setup, $B = 256$ and $B_{\text{mini}} = 64$,
 428 so each train-batch step contains $M = B/B_{\text{mini}} = 4$ optimizer updates. This conversion is useful
 429 for implementation accounting, but it is not the staleness axis used in the plots.

430 We use k as the x-axis because it is the controlled train-batch staleness intervention shared across
 431 methods. The sweep covers $k \in \{0, 1, 2, 4, 8, 16, 32, 64, 128\}$ across the forward-KL, reverse-KL /
 432 PPO-style, M2PO / DecPPO, top- k , and Monte-Carlo support-size variants; apart from the estimator
 433 choice and k , these runs share the common model, data, batch-size, generation, and evaluation
 434 settings in Table 3.

435 C Datasets and Metrics

436 **Training data.** We filter the DeepMath dataset [8] to retain 57,630 math problems with difficulty
 437 level greater than or equal to 6, and use this filtered subset as the training data.

438 **Evaluation datasets.** Table 4 lists the evaluation datasets used: AIME 2024 [30], AIME 2025 [31],
 439 and AMC 2023 [13]. AIME24 is evaluated every 20 steps. The remaining datasets are only evaluated
 440 for the final checkpoint.

Table 5: Existing assets used in this work, with source URLs, license names, and versions.

Category	Asset	URL	License name	Version
Dataset	DeepMath	https://huggingface.co/datasets/zwe99/DeepMath-103K	MIT	2025
Dataset	AIME 2024	https://huggingface.co/datasets/Maxwell-Jia/AIME_2024	MIT	2024
Dataset	AIME 2025	https://huggingface.co/datasets/yentinglin/aime_2025	Not specified	2025
Dataset	AMC 2023 / American Mathematics Competitions	https://huggingface.co/datasets/math-ai/amc23	Not specified	2023
Model	Qwen3-4B-Base	https://huggingface.co/Qwen/Qwen3-4B-Base	Apache-2.0	Not specified
Model	Qwen3-1.7B-Base	https://huggingface.co/Qwen/Qwen3-1.7B-Base	Apache-2.0	Not specified
Model	Qwen3-8B-Base	https://huggingface.co/Qwen/Qwen3-8B-Base	Apache-2.0	Not specified
Model	Qwen3-4B	https://huggingface.co/Qwen/Qwen3-4B	Apache-2.0	Not specified
Model	Qwen3-1.7B	https://huggingface.co/Qwen/Qwen3-1.7B	Apache-2.0	Not specified
Model	Qwen3-8B	https://huggingface.co/Qwen/Qwen3-8B	Apache-2.0	Not specified
Model	Qwen3-30B-A3B-Instruct-2507	https://huggingface.co/Qwen/Qwen3-30B-A3B-Instruct-2507	Apache-2.0	Qwen3-2507
Software	vLLM	https://github.com/vllm-project/vllm	Apache-2.0	0.16.0
Software	PyTorch / PyTorch FSDP	https://github.com/pytorch/pytorch	BSD-3-Clause	2.9.1

441 **Accuracy metric.** Evaluation samples 32 responses per problem. For a dataset D , the reported
 442 Avg@32 is the mean per-problem pass rate,

$$\text{Avg@32}(D) = 100 \cdot \frac{1}{|D|} \sum_{i \in D} \frac{c_i}{32}, \quad (14)$$

443 where c_i is the number of sampled responses judged correct for problem i . Paper tables and plots use
 444 Avg@32 unless noted otherwise.

445 D Existing Asset Licenses

446 Table 5 lists the reused assets.

447 E Importance-Sampling Ablation

448 Before treating multi-sampling as an improvement, we separate it from IS. Increasing m changes the
 449 Monte Carlo variance of the estimator. It does not define the target distribution. Old-to-current IS is
 450 still the mechanism that turns behavior-policy samples into an estimator of the current reverse-KL
 451 local gradient. Figure 8 compares MC1 and MC16 with and without IS to test this distinction directly.

452 F AsyncOPD Scheduler Details

453 This appendix gives the implementation details omitted from Section 7. Our scheduler, AsyncOPD,
 454 follows the fully asynchronous systems structure of AReal [4], but the queue contains OPD cache
 455 items rather than reward-labeled RL trajectories.

456 **Queue interface.** The pipeline has three long-running stages: rollout generation, teacher scoring,
 457 and learner training. Rollout workers sample trajectories from their latest synchronized student
 458 snapshot. For each visited prefix s , they cache the MC actions, rollout log probabilities under p_{old} ,
 459 and the rollout student version. The main scheduler comparison uses MC64; the MC1 runs use the
 460 same queue interface with one cached action. The teacher scores the cached actions. The learner

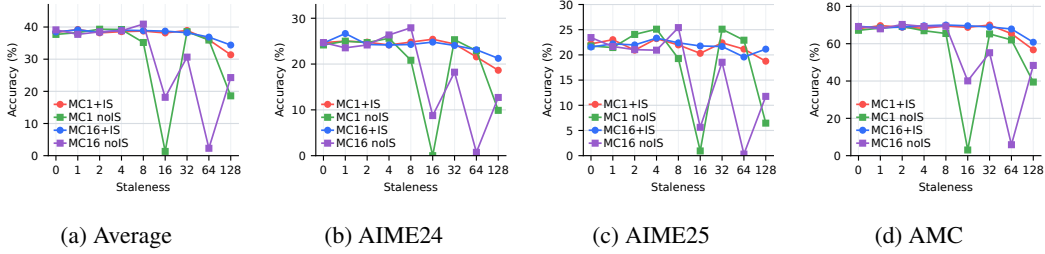


Figure 8: Accuracy comparison under staleness for MC importance-sampling ablations. Increasing the number of samples reduces Monte Carlo variance, but old-to-current IS is still needed to correct stale-policy sampling.

461 then recomputes $\log p_\theta(a | s)$ and $A_\theta(a, s)$ under the current student, and applies the unclipped
 462 old-to-current IS estimator from Section 6.

463 **Weight synchronization and queue capacity.** During AsyncOPD weight synchronization, rollout
 464 workers pause generation. In the keep-mode used for the scheduler experiments, in-flight requests
 465 are not discarded: the already sampled token prefix is kept, the student weights are updated, and
 466 the running-request prefix cache is reset so the engine rebuilds the attention state for that prefix
 467 under the new weights before generation resumes. Thus, tokens before the synchronization point
 468 are reused rather than regenerated, while later tokens are sampled under the new student snapshot.
 469 Each completed sample records the token index at which the active weight version changes. The
 470 queue-depth parameter τ is enforced as a capacity bound rather than as a learner-side drop rule.
 471 The coordinator creates a semaphore with $(\tau + 1)B$ permits, where B is the effective train batch
 472 size. The prompt feeder acquires one permit before submitting a prompt to rollout, and the train
 473 dispatcher releases permits only after the corresponding samples have been consumed by a learner
 474 update. During weight synchronization, a sync gate prevents the feeder from using newly released
 475 permits until rollout workers have received the updated weights. Thus, smaller τ limits the amount of
 476 unconsumed rollout work in the pipeline, while larger τ permits a deeper backlog and more overlap.
 477 The queues themselves remain FIFO; items are not evicted for being stale.

478 **Training throughput metric.** The table reports training throughput. Let n_j be the number of
 479 response tokens used by learner update j , and let t_j be the train wall-clock time after that update.
 480 Discarding the first five warmup updates, we compute

$$\text{throughput} = \frac{\sum_{j=6}^J n_j}{\sum_{j=6}^J (t_j - t_{j-1})} = \frac{\sum_{j=6}^J n_j}{t_J - t_5}.$$

481 Speedups are normalized to the strict-sync run with the same student and MC setting.

482 **Pipeline overlap metric.** Let $\mathcal{S} = \{\text{rollout}, \text{teacher}, \text{train}\}$. For teacher and train, we merge
 483 overlapping busy intervals within each stage and compute the merged busy time T_s . Rollout has N_r
 484 workers, so we first merge intervals within each worker i and then define the rollout-stage busy time
 485 as the worker-normalized average

$$T_{\text{rollout}} = \frac{1}{N_r} \sum_{i=1}^{N_r} T_{\text{rollout},i}.$$

486 Let T_{wall} be the elapsed train wall-clock interval from the first to the last recorded pipeline-stage
 487 interval. We define

$$\text{overlap} = \frac{\sum_{s \in \mathcal{S}} T_s}{T_{\text{wall}}}.$$

488 A mostly serial schedule has overlap near 1, and the maximum remains 3: all rollout workers, teacher
 489 scoring, and training busy for the full interval.

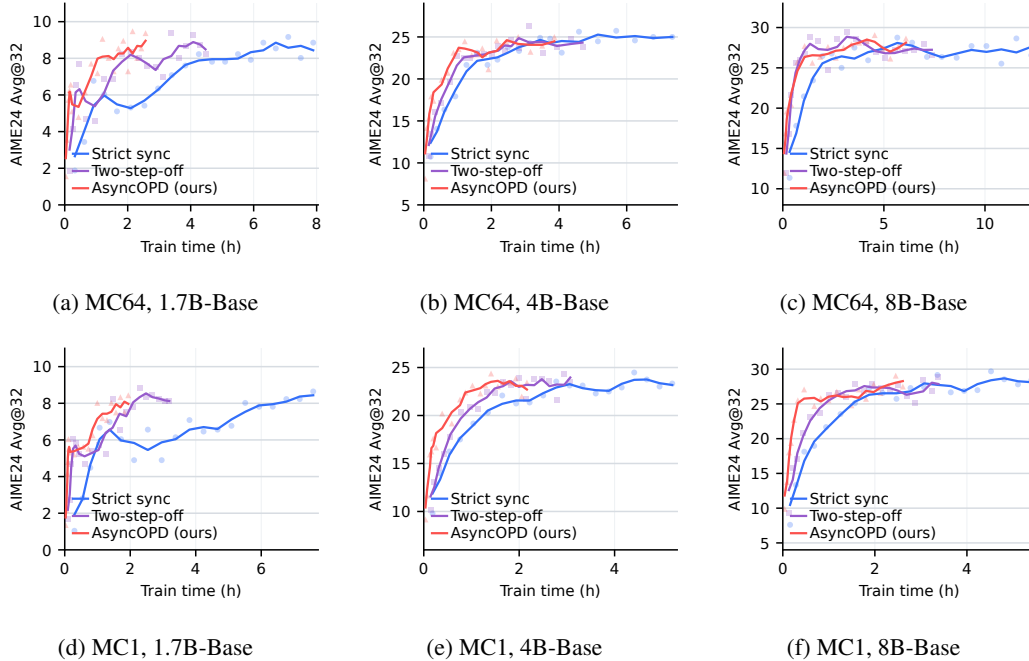


Figure 9: Train-time AIME24 Avg@32 for Qwen3-Base students with MC64 and MC1. Lines are 3-point moving averages; faint markers are raw evaluations; colors denote scheduler. AsyncOPD reaches later checkpoints sooner, so its accuracy improves earlier in wall-clock time.

490 **Hardware and testing protocol.** Each scheduler run uses one 8-GPU node. One GPU is reserved
 491 for teacher scoring. The remaining seven GPUs are the rollout/training pool. Rollout generation uses
 492 data parallelism, and learner training uses PyTorch FSDP. Strict sync runs time-share this pool: all
 493 seven GPUs run rollout, then all seven switch to training, and the cycle repeats. The two-step-off
 494 and our AsyncOPD runs split the same seven GPUs concurrently: 4 GPUs for rollout workers and 3
 495 GPUs for the FSDP trainer.

496 For each student size and MC setting, we compare strict sync, two-step-off, and our AsyncOPD
 497 scheduler with the same teacher, training data, evaluation metrics, and reverse-KL estimator: current-
 498 policy A_θ , no clipping, and old-to-current IS correction. Two-step-off fixes a two-update offset
 499 between rollout and the learner update that consumes it, so stale rollout reuse is static and controlled
 500 rather than produced by queue timing. We use this offset because it is the fastest static step-off
 501 schedule under the 4-rollout/3-trainer split. The OPD pipeline has three serial stages: rollout
 502 generation, teacher scoring, and learner training. Therefore, a two-step offset is enough to keep all
 503 stages occupied in the gated schedule. Larger offsets only make the consumed data older; they do not
 504 create another OPD stage to overlap or remove the step-off batch barrier. We measure final-checkpoint
 505 Avg@32 and train wall-clock time over the same training horizon.

506 **Qwen3-Base train-time accuracy.** Figure 9 provides the train-time view for Qwen3-Base students.
 507 AsyncOPD reaches later checkpoints sooner, so accuracy improves earlier in wall-clock time across
 508 student sizes and MC settings.

509 **Additional Qwen3 AsyncOPD results.** For the Qwen3 1.7B, 4B, and 8B student rows, we disable
 510 thinking at the tokenizer prompt-formatting level: prompt construction uses the Qwen3 tokenizer’s
 511 non-thinking chat-template mode before rollout and evaluation. Table 6 and Fig. 10 report this
 512 comparison. The systems pattern matches the Qwen3-Base results: AsyncOPD has the highest
 513 throughput and overlap, reaching up to $3.8\times$ strict-sync throughput on MC64 and up to $3.2\times$ on
 514 MC1. The train-time accuracy plots show the same wall-clock pattern as the main Qwen3-Base
 515 results: AsyncOPD reaches later checkpoints sooner, so accuracy improves earlier across student
 516 sizes and MC settings.

Table 6: Additional AsyncOPD scheduler results for Qwen3 students with thinking disabled. Train tok/s is training throughput; parentheses show speedup over the matched strict-sync baseline. Overlap is concurrent OPD-stage activity. Avg@32 is final AIME24. AsyncOPD achieves the highest throughput and overlap in all matched settings while maintaining comparable final accuracy.

Student	Scheduler	MC64			MC1		
		Train tok/s (\times sync)	Overlap	Avg@32	Train tok/s (\times sync)	Overlap	Avg@32
1.7B	Strict sync	11.2k (1.00 \times)	0.78	35.00	11.7k (1.00 \times)	0.78	35.21
	Two-step-off	18.0k (1.61 \times)	1.58	32.81	20.8k (1.78 \times)	1.71	33.54
	AsyncOPD (ours)	24.2k (2.16\times)	2.07	33.23	37.0k (3.17\times)	2.56	34.79
4B	Strict sync	4.5k (1.00 \times)	0.87	54.90	7.4k (1.00 \times)	0.83	56.15
	Two-step-off	11.5k (2.58 \times)	1.54	56.77	13.0k (1.75 \times)	1.63	58.02
	AsyncOPD (ours)	17.0k (3.82\times)	2.15	54.69	18.5k (2.49\times)	2.26	56.25
8B	Strict sync	5.6k (1.00 \times)	0.86	59.69	3.9k (1.00 \times)	0.87	58.96
	Two-step-off	7.5k (1.34 \times)	1.45	58.33	9.6k (2.45 \times)	1.61	60.73
	AsyncOPD (ours)	10.1k (1.81\times)	2.11	60.52	10.9k (2.78\times)	2.17	58.65

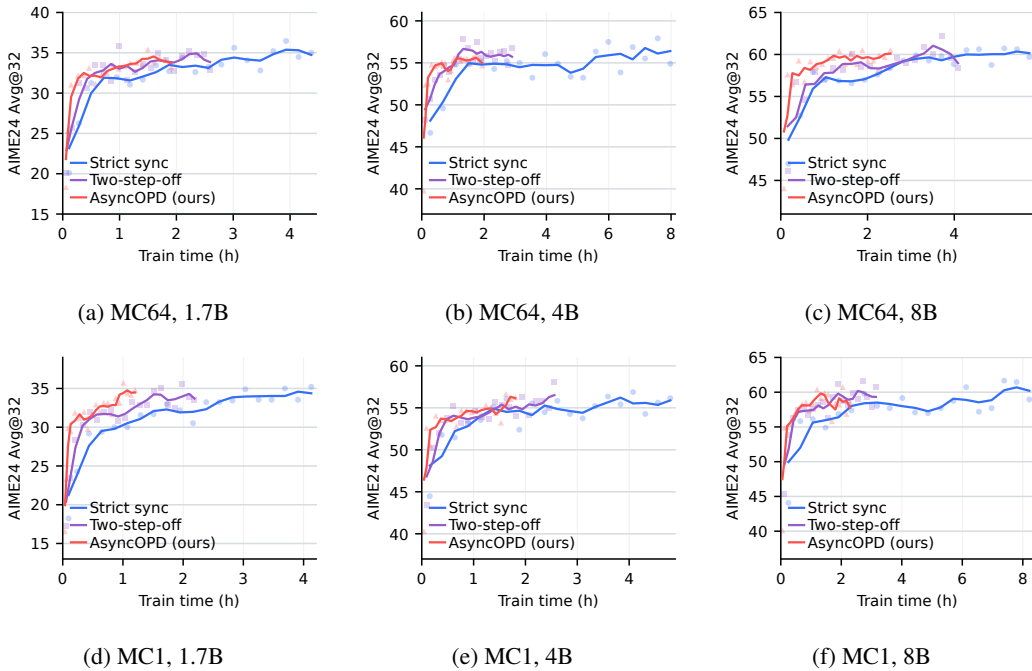


Figure 10: Train-time AIME24 Avg@32 for Qwen3 1.7B, 4B, and 8B students with thinking disabled, using MC64 and MC1. Lines are 3-point moving averages; faint markers are raw evaluations; colors denote scheduler. AsyncOPD reaches later checkpoints sooner, so its accuracy improves earlier in wall-clock time.

517 **NeurIPS Paper Checklist**

518 **1. Claims**

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

521 Answer: [Yes]

522 Justification: The abstract and introduction accurately reflect the paper’s main contributions
523 and scope: the staleness analysis of asynchronous OPD, the forward-/reverse-KL compari-
524 son, the reverse-KL surrogate study, multi-sample MC, and the AsyncOPD system. These
525 claims are supported by the main experiments and analysis in Sections 4–7.

526 Guidelines:

- 527 • The answer [N/A] means that the abstract and introduction do not include the claims
528 made in the paper.
- 529 • The abstract and/or introduction should clearly state the claims made, including the
530 contributions made in the paper and important assumptions and limitations. A [No] or
531 [N/A] answer to this question will not be perceived well by the reviewers.
- 532 • The claims made should match theoretical and experimental results, and reflect how
533 much the results can be expected to generalize to other settings.
- 534 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
535 are not attained by the paper.

536 **2. Limitations**

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

538 Answer: [Yes]

539 Justification: The paper discusses limitations and future work in Section 8, including
540 the focus on sparse and Monte Carlo OPD estimators and the absence of full-vocabulary
541 dense-KL evaluation in the asynchronous pipeline.

542 Guidelines:

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

569 **3. Theory assumptions and proofs**

570 Question: For each theoretical result, does the paper provide the full set of assumptions and
571 a complete (and correct) proof?

572 Answer: [Yes]

573 Justification: For the theoretical and estimator-level claims, the paper states the relevant OPD
574 setup and assumptions and provides derivations for the forward-/reverse-KL formulations,
575 the reverse-KL policy-gradient form, and the old-to-current importance-sampling identity.
576 The paper does not rely on separate theorem-style results beyond these analytical derivations.

577 Guidelines:

- 578 • The answer [N/A] means that the paper does not include theoretical results.
- 579 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
580 referenced.
- 581 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 582 • The proofs can either appear in the main paper or the supplemental material, but if
583 they appear in the supplemental material, the authors are encouraged to provide a short
584 proof sketch to provide intuition.
- 585 • Inversely, any informal proof provided in the core of the paper should be complemented
586 by formal proofs provided in appendix or supplemental material.
- 587 • Theorems and Lemmas that the proof relies upon should be properly referenced.

588 4. Experimental result reproducibility

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

592 Answer: [Yes]

593 Justification: The paper discloses the main information needed to reproduce the experimental
594 results, including model choices, datasets, training horizon, optimization settings, generation
595 settings, evaluation protocol, implementation stack, scheduler configuration, and hardware
596 setup. Additional estimator and scheduler details are provided in the appendices.

597 Guidelines:

- 598 • The answer [N/A] means that the paper does not include experiments.
- 599 • If the paper includes experiments, a [No] answer to this question will not be perceived
600 well by the reviewers: Making the paper reproducible is important, regardless of
601 whether the code and data are provided or not.
- 602 • If the contribution is a dataset and/or model, the authors should describe the steps taken
603 to make their results reproducible or verifiable.
- 604 • Depending on the contribution, reproducibility can be accomplished in various ways.
605 For example, if the contribution is a novel architecture, describing the architecture fully
606 might suffice, or if the contribution is a specific model and empirical evaluation, it may
607 be necessary to either make it possible for others to replicate the model with the same
608 dataset, or provide access to the model. In general, releasing code and data is often
609 one good way to accomplish this, but reproducibility can also be provided via detailed
610 instructions for how to replicate the results, access to a hosted model (e.g., in the case
611 of a large language model), releasing of a model checkpoint, or other means that are
612 appropriate to the research performed.
- 613 • While NeurIPS does not require releasing code, the conference does require all submis-
614 sions to provide some reasonable avenue for reproducibility, which may depend on the
615 nature of the contribution. For example
 - 616 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
617 to reproduce that algorithm.
 - 618 (b) If the contribution is primarily a new model architecture, the paper should describe
619 the architecture clearly and fully.
 - 620 (c) If the contribution is a new model (e.g., a large language model), then there should
621 either be a way to access this model for reproducing the results or a way to reproduce
622 the model (e.g., with an open-source dataset or instructions for how to construct
623 the dataset).

624 (d) We recognize that reproducibility may be tricky in some cases, in which case
625 authors are welcome to describe the particular way they provide for reproducibility.
626 In the case of closed-source models, it may be that access to the model is limited in
627 some way (e.g., to registered users), but it should be possible for other researchers
628 to have some path to reproducing or verifying the results.

629 5. Open access to data and code

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

633 Answer: [Yes]

634 Justification: The paper states that the AsyncOPD implementation is released through an
635 anonymized repository, and the datasets and models used in the experiments are public
636 or cited. The appendix provides the experimental setup needed to reproduce the reported
637 results.

638 Guidelines:

- 639 • The answer [N/A] means that paper does not include experiments requiring code.
- 640 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
641 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 642 • While we encourage the release of code and data, we understand that this might not
643 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
644 including code, unless this is central to the contribution (e.g., for a new open-source
645 benchmark).
- 646 • The instructions should contain the exact command and environment needed to run to
647 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 648 • The authors should provide instructions on data access and preparation, including how
649 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 650 • The authors should provide scripts to reproduce all experimental results for the new
651 proposed method and baselines. If only a subset of experiments are reproducible, they
652 should state which ones are omitted from the script and why.
- 653 • At submission time, to preserve anonymity, the authors should release anonymized
654 versions (if applicable).
- 655 • Providing as much information as possible in supplemental material (appended to the
656 paper) is recommended, but including URLs to data and code is permitted.

658 6. Experimental setting/details

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

661 Answer: [Yes]

662 Justification: The paper specifies the training and test settings, including the filtered train-
663 ing data, prompt and response lengths, optimizer, rollout generation settings, evaluation
664 sampling protocol, implementation stack, and hardware allocation (see Appendix and code
665 implementation).

666 Guidelines:

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

672 7. Experiment statistical significance

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

675 Answer: [No]

676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728

Justification: The main experiments are computationally expensive, so repeating every experiment across multiple random seeds was not feasible within our available compute budget. We acknowledge this limitation.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

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

Answer: [Yes]

Justification: The paper provides the compute setup used for the experiments, including the single 8×B200 node, GPU allocation across teacher scoring, rollout, and training, and the throughput and overlap measurements used in the scheduler comparisons (see Appendix).

Guidelines:

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

9. Code of ethics

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

Answer: [Yes]

Justification: The research follows standard experimental practice for LLM post-training, uses public or cited datasets and models, preserves submission anonymity, and does not involve human-subject experiments or unsafe data collection.

Guidelines:

- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.

- 729 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
730 eration due to laws or regulations in their jurisdiction).

731 **10. Broader impacts**

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

734 Answer: [N/A]

735 Justification: The paper discusses positive impact through improved training efficiency
736 and hardware utilization for LLM post-training, and the work is methodological and
737 systems-focused rather than a deployment of a new end-user model or application. No
738 new application-specific societal risks are introduced beyond those of the underlying LLM
739 post-training setting.

740 Guidelines:

- 741 • The answer [N/A] means that there is no societal impact of the work performed.
- 742 • If the authors answer [N/A] or [No], they should explain why their work has no societal
743 impact or why the paper does not address societal impact.
- 744 • Examples of negative societal impacts include potential malicious or unintended uses
745 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
746 (e.g., deployment of technologies that could make decisions that unfairly impact specific
747 groups), privacy considerations, and security considerations.
- 748 • The conference expects that many papers will be foundational research and not tied
749 to particular applications, let alone deployments. However, if there is a direct path to
750 any negative applications, the authors should point it out. For example, it is legitimate
751 to point out that an improvement in the quality of generative models could be used to
752 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
753 that a generic algorithm for optimizing neural networks could enable people to train
754 models that generate Deepfakes faster.
- 755 • The authors should consider possible harms that could arise when the technology is
756 being used as intended and functioning correctly, harms that could arise when the
757 technology is being used as intended but gives incorrect results, and harms following
758 from (intentional or unintentional) misuse of the technology.
- 759 • If there are negative societal impacts, the authors could also discuss possible mitigation
760 strategies (e.g., gated release of models, providing defenses in addition to attacks,
761 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
762 feedback over time, improving the efficiency and accessibility of ML).

763 **11. Safeguards**

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

767 Answer: [N/A]

768 Justification: The paper does not release a high-risk generative model, scraped dataset, or
769 other asset requiring special safeguards for responsible release. The released asset is a
770 training-system implementation.

771 Guidelines:

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

782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833

12. Licenses for existing assets

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

Answer: [Yes]

Justification: We cite the original creators and sources of the existing assets used in this work. Table 5 lists the source URL, license name, and version for each reused asset.

Guidelines:

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

13. New assets

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

Answer: [Yes]

Justification: The paper introduces and releases AsyncOPD, a fully asynchronous OPD training pipeline. The implementation, estimator choices, scheduler behavior, and experimental setup are documented in the main text and appendices. Instructions on how to reproduce the results are provided in the code.

Guidelines:

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

14. Crowdsourcing and research with human subjects

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

Answer: [N/A]

Justification: The paper does not involve crowdsourcing or research with human subjects.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.

834 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
835 or other labor should be paid at least the minimum wage in the country of the data
836 collector.

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

839 Question: Does the paper describe potential risks incurred by study participants, whether
840 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
841 approvals (or an equivalent approval/review based on the requirements of your country or
842 institution) were obtained?

843 Answer: [N/A]

844 Justification: The paper does not involve crowdsourcing or human-subject research, so IRB
845 approval or equivalent review is not applicable.

846 Guidelines:

- 847 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
848 with human subjects.
- 849 • Depending on the country in which research is conducted, IRB approval (or equivalent)
850 may be required for any human subjects research. If you obtained IRB approval, you
851 should clearly state this in the paper.
- 852 • We recognize that the procedures for this may vary significantly between institutions
853 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
854 guidelines for their institution.
- 855 • For initial submissions, do not include any information that would break anonymity (if
856 applicable), such as the institution conducting the review.

857 **16. Declaration of LLM usage**

858 Question: Does the paper describe the usage of LLMs if it is an important, original, or
859 non-standard component of the core methods in this research? Note that if the LLM is used
860 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
861 scientific rigor, or originality of the research, declaration is not required.

862 Answer: [N/A]

863 Justification: LLMs were not used as an important, original, or non-standard component for
864 developing the research methodology or manuscript. LLMs were used only for writing, or
865 editing (declaration is not required under the checklist policy).

866 Guidelines:

- 867 • The answer [N/A] means that the core method development in this research does not
868 involve LLMs as any important, original, or non-standard components.
- 869 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
870 be described.