
Good Experience Maximization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Post-training language models typically follows one of two paradigms: reinforcement
2 learning (RL) and supervised fine-tuning (SFT). Although typically thought
3 of as different camps, we show that both paradigms can be viewed as target distri-
4 bution matching. We use this perspective to analyze the tension between policy
5 gradient and buffer-based SFT and establish that a good post-training target should
6 preserve past successes while continuing to incorporate newly-discovered ones.
7 We introduce Good Experience Maximization (GEM), a simple framework that
8 resolves this tension by training the policy toward a mixture of historical and
9 newly generated successes. This construction is amenable to a variety of simple
10 algorithmic choices, is agnostic to the sampling process that populates the buffer
11 and naturally subsumes several existing post-training methods as special cases. We
12 prove that the GEM target yields globally optimal expected reward at convergence
13 under realizability, and we derive a flow-based update that guarantees one-step
14 policy improvement. This is a distinctively different update to SFT and relies on
15 constructing a straight line “path” toward the target distribution. Empirically both
16 the “flow” and the cross entropy version of GEM perform on par with leading
17 post-training methods and out-perform it in sparse reward settings.

18 1 Introduction

19 Good experiences are worth repeating. The brain seems to operate on this principle: after salient
20 experiences, it spontaneously reactivates these memories offline, prioritizing those tied to reward,
21 novelty, or uncertainty [Mattar and Daw, 2018]. This replay has many benefits; recent evidence
22 shows it supports important cognitive tasks, including memory consolidation, planning, and decision-
23 making [Momennejad et al., 2018, Schapiro et al., 2018]. Remarkably, this mechanism enables a
24 single meaningful experience to continue shaping future behavior without requiring us to physically
25 reenact it. Machine learning has a direct analog of this idea: *experience replay* [Lin, 1992, Mnih et al.,
26 2013]. Rather than being discarded after a single use, past interactions are stored in a replay buffer
27 and reused during training, which enables experiences to inform many gradient updates and amortizes
28 the cost of data collection. This simple mechanism has supported major empirical successes across
29 challenging decision-making domains [Andrychowicz et al., 2017, Vinyals et al., 2019].

30 Despite these advantages, current methods for post-training language models (LMs) are dominated
31 by methods that require online samples at every step [Williams, 1992, Shao et al., 2024, Tajwar et al.,
32 2026]. These on-policy methods generate K fresh completions, or rollouts, per problem x from the
33 current LM (or policy) π_θ , compute a gradient update, and discard the data. However, this process is
34 wasteful: a correct reasoning trace discovered at a particular training step may remain a useful signal
35 long after the LM has been updated. To continue improving, the model must rediscover solutions it
36 has already found, which is particularly problematic when solutions are rare. Furthermore, discarding
37 these solutions risks *forgetting*: if the model is no longer trained on traces that solved a given problem,
38 it may lose the ability to solve it [Yue et al., 2025]. An ideal training regime would resemble human
39 memory: hold onto what worked, and keep learning from it.

40 Reuse alone has its own pathology: a learner that solely revisits past successes risks never discovering
 41 new ones. Indeed, a learner trained only on historical successes is constrained by the coverage of the
 42 data and performance can degrade if the policy moves beyond that distribution [Ross et al., 2011]. But
 43 relying on the policy itself as the only memory of past successes is also insufficient. If a solution is
 44 not externally stored, then recovering it requires sampling it again from the model, which is unlikely
 45 for rare reasoning traces. In principle, one may try to preserve past behavior through regularization
 46 or value function learning, but extracting solutions with them still requires search or sampling. An
 47 effective post-training regime therefore requires both memory and exploration: the LM ought to
 48 continue learning from found successes while still producing new ones worth remembering.

49 **Contributions.** This paper studies how to realize this principle in post-training. We introduce
 50 *Good Experience Maximization* (GEM), a framework that incorporates memory through a replay
 51 buffer of past successes and exploration by collecting online data at every step. We first formalize
 52 post-training as target distribution matching and use this lens to study the failure cases of policy
 53 gradient and buffer-based SFT (Section 3). We then use this lens to develop GEM, a simple post-
 54 training framework that trains the policy on a mixture of historical and newly generated successes
 55 (Section 4). We prove that this mixture target yields global expected reward at convergence and
 56 develop a novel flow-based update that establishes improvement at each update step. Empirically, we
 57 provide a practical algorithm and show that GEM can match or exceed RL algorithms in terms of
 58 performance and sample efficiency (Section 5) In Section 6, we demonstrate how this framework
 59 subsumes a broad set of existing post-training methods as special cases, including MaxRL [Tajwar
 60 et al., 2026] and Iterated STaR [Zelikman et al., 2022].

61 2 Background and Preliminaries

62 **Problem setup.** We study correctness-based problems, which can be abstracted as binary successes
 63 or failures for each input \mathbf{x} . This setting is known as reinforcement learning with verifiable rewards
 64 (RLVR) [Lambert et al., 2024, Su et al., 2025] and encompasses a wide set of modern problems. In
 65 this setting, \mathcal{X}, \mathcal{Y} denotes the input and output spaces, respectively, and we are given query access to
 66 a ground-truth reward function $R : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$. The goal is to maximize the expected reward of
 67 some parameterized policy $\pi_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ given some distribution $\mathbf{x} \sim \mu$,

$$\mathcal{J}(\pi_\theta) = \mathbb{E}_{\mathbf{x} \sim \mu, \mathbf{y} \sim \pi_\theta(\cdot | \mathbf{x})} [R(\mathbf{x}, \mathbf{y})].$$

68 In many modern settings, the model does not sample directly from \mathcal{Y} ; instead, it may generate a trace z ,
 69 which is deterministically parsed into a reasoning trace \mathbf{r} and \mathbf{y} . We can treat each \mathbf{r} as a latent variable
 70 that is marginalized out to obtain the answer distribution $\pi_\theta(\mathbf{y} | \mathbf{x}) = \sum_{\mathbf{r}} \pi_\theta(\mathbf{r} | \mathbf{x}) \pi_\theta(\mathbf{y} | \mathbf{x}, \mathbf{r})$.

71 **Policy optimization.** There are two dominant paradigms for post-training generative models. One
 72 of these is online policy-gradient methods [Williams, 1992, Schulman et al., 2017, Guo et al., 2025],
 73 which estimate gradients of the policy optimization objective as

$$\nabla_\theta \mathcal{J}(\theta) = \mathbb{E}_{\mathbf{x}, \mathbf{y}} [A(\mathbf{x}, \mathbf{y}) \nabla_\theta \log \pi_\theta(\mathbf{y} | \mathbf{x})], \quad (1)$$

74 where $A(\mathbf{x}, \mathbf{y})$ denotes some advantage estimate. Another, related perspective instead maximizes a
 75 lower bound on the probability of success [Dayan and Hinton, 1997, Abdolmaleki et al., 2018a,b,
 76 2024, Tajwar et al., 2026]. The other paradigm follows the simple principle of “sample, filter, then
 77 SFT”. In practice, methods typically sample many answers across a large dataset, filter for the set
 78 of good answers, then maximize the likelihood of this filtered dataset [Brandfonbrener et al., 2021,
 79 Zelikman et al., 2022, Dong et al., 2023, Gulcehre et al., 2023, Singh et al., 2023, Havrilla et al.,
 80 2024]. Given a buffer,

$$\mathcal{B} := \{(\mathbf{x}, \mathbf{y}) : \mathbf{y} \sim \pi(\cdot | \mathbf{x}), R^*(\mathbf{x}, \mathbf{y}) = 1\},$$

81 the model is trained with the supervised objective:

$$\mathcal{L}_{\text{SFT}}(\theta) := -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P_{\mathcal{B}}} [\log \pi_\theta(\mathbf{y} | \mathbf{x})], \quad (2)$$

82 with the gradient:

$$\nabla_\theta \mathcal{L}_{\text{SFT}}(\theta) := -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P_{\mathcal{B}}} [\nabla_\theta \log \pi_\theta(\mathbf{y} | \mathbf{x})]. \quad (3)$$

83 Equations (1) and (3) provide the objects we analyze in this work. Although policy gradient and
 84 filtered SFT differ operationally, both gradients are expectations of $\nabla_\theta \log \pi_\theta(\mathbf{y} | \mathbf{x})$ under some
 85 distribution over prompt-answer pairs. We use this shared form in Section 6 to study post-training as
 86 a problem of target distribution design.

87 3 A Target-Matching View

88 While the community has explored a wide range of policy optimization approaches (Section 2), we
 89 observe that many can be viewed as iterative update procedures that move a policy toward a target
 90 distribution. We refer to this formulation colloquially as the ‘‘SFT update’’. Specifically, in the binary
 91 reward case, regular policy gradient [Williams, 1992], expectation-maximization based methods that
 92 maximize the likelihood of success [Dayan, 1993, Abdolmaleki et al., 2018b, Tajwar et al., 2026],
 93 as well as filtered SFT methods [Brandfonbrener et al., 2021, Dong et al., 2023, Gulcehre et al.,
 94 2023, Havrilla et al., 2024] can all be viewed as minimizing a per-iteration SFT loss toward the target
 95 distribution at each input \mathbf{x} ,

$$q_t(\mathbf{y}|\mathbf{x}) = \frac{\pi_t(\mathbf{y}|\mathbf{x})R(\mathbf{x}, \mathbf{y})}{\mathbb{E}_{\pi_t(\mathbf{y}|\mathbf{x})}[R(\mathbf{x}, \mathbf{y})]}, \quad (4)$$

96 We expand on this connection in Section 6. Intuitively, these methods all follow a gradient of the
 97 log probability, and we can absorb any scaling to the per-sample log probability into the probability
 98 of the sampling distribution. This target-matching view exposes a tension between two desirable
 99 properties, which we will now discuss.

100 3.1 Policy Gradient Creates Good Experiences But Does Not (Explicitly) Retain Them

101 Policy gradient samples small batches online and typically
 102 perform single updates. A successful trace sampled at
 103 iteration t affects the policy only through the update taken
 104 at that iteration. If the trace is not externally stored, then
 105 for it to affect future updates, the policy must sample it
 106 again. Thus, the probability that \mathbf{y} contributes repeatedly
 107 to learning remains proportional to the mass the current
 108 policy places on it $\pi_t(\mathbf{y}|\mathbf{x})$. Thus, for any good solutions
 109 to be repeatedly reinforced, the policy must repeatedly
 110 sample them, making learning slow and sample-inefficient.

111 Yet, the single update is also the benefit of policy gradi-
 112 ent algorithms. Because each update is based on fresh
 113 samples from the current policy, policy gradient continues
 114 enjoys full support over problem distribution. Moreover,
 115 a successful generation can produce an immediate update:
 116 in the binary reward setting, even a single high-reward
 117 sample can increase the likelihood of the behavior that
 118 produced it. Thus, learning can begin before the learner
 119 has collected successful outputs for every input. With
 120 function approximation, this update may also generalize
 121 across related inputs, allowing progress on questions that
 122 have not yet been solved directly.

123 3.2 Iterative SFT Retains Good 124 Experiences but Does Not Guarantee Improvement

125 Alternatively, the canonical class of ‘‘filtered SFT’’ algo-
 126 rithms [Zelikman et al., 2022, Brandfonbrener et al., 2021,
 127 Dong et al., 2023, Gulcehre et al., 2023, Havrilla et al.,
 128 2024] typically sample (nearly) the full dataset, then do
 129 multiple gradient steps on the filtered dataset. As we estab-
 130 lish above and in Section 6, a single step of this procedure
 131 is in expectation *equivalent* to a single step of online policy gradient (or an online variant such as
 132 EM or MaxRL, depending on if we sample by question first, or by correct answers). The distinct
 133 advantage of iterative SFT is the ability to reuse successful generations across multiple gradient
 134 updates. Filtered SFT therefore has two appealing properties: in theory, it can (i) perform multiple
 135 updates, and (ii) make use of multiple sources of ‘‘good’’ data. But when does fitting an accumulated
 136 buffer actually improve the policy?

137 One limitation is coverage. If we train to *exactly* match only the buffer distribution $P_{\mathcal{B}}$, we are limited
 138 by the support of \mathcal{B} . We can write down a simple bound for the amount of reward we can guarantee
 139 *on the training set* \mathcal{D} where \mathcal{B} provides no signal.

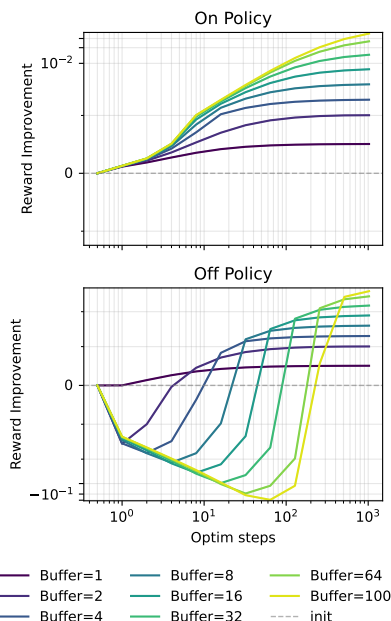


Figure 1: SFT training dynamics in a toy categorical task. Figures illustrate SFT updates toward a buffer distribution with (i) high reward samples filtered from the current policy, or (ii) high reward samples from a very different policy. Colors indicate the state space coverage of the buffer.

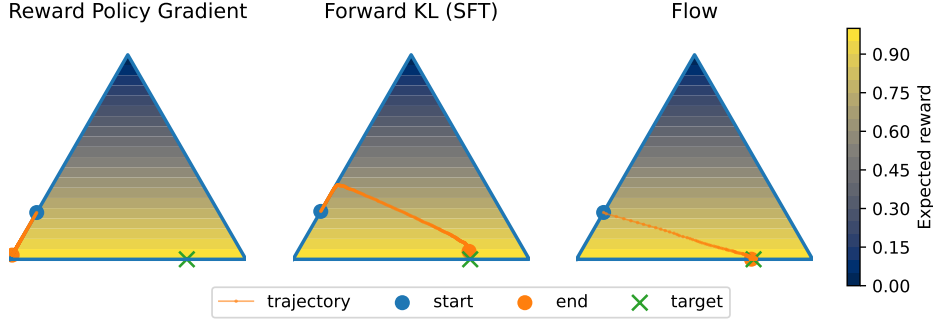


Figure 2: Geometry of optimization with on and off policy data on a 3 categories simplex. Example of reward policy gradient, forward KL (SFT) toward a target distribution, and our flow update.

140 *Remark 3.1.* Suppose we have binary rewards, training dataset \mathcal{D} , and a buffer distribution $P_{\mathcal{B}}$ whose
 141 support is limited to $M < |\mathcal{D}|$ questions. If a policy π^* perfectly approximates $P_{\mathcal{B}}$, the minimum
 142 expected reward it achieves $\mathcal{J}(\pi^*) = \frac{M}{|\mathcal{D}|}$. For the remaining $|\mathcal{D}| - M$ questions, the objective is
 143 undefined, and $\pi^*(\cdot|\mathbf{x})$ may behave arbitrarily poorly.

144 The above is reflected in Figure 1: the *eventual* performance grows monotonically with the coverage
 145 of the buffer over the question set. Thus, with a sufficiently broad buffer and sufficient training,
 146 filtered SFT can achieve high training-set reward. The difficulty is that obtaining such coverage may
 147 require sampling nearly the full dataset before substantial learning can occur.

148 Another limitation is monotonic improvement. To avoid repeatedly sampling every input, we may
 149 want to reuse data from older policies or other sources. As shown in Figure 1, if the buffer data
 150 is good but “off policy” (in the sense that it has low likelihood under the current policy), we can
 151 end up with optimization dynamics where in a small number of steps we will end up with *lower*
 152 *rewards*, despite the buffer having higher expected rewards. The issue is not that the buffered samples
 153 are bad; rather, a standard SFT path toward an arbitrary buffer distribution is not necessarily an
 154 improvement-preserving path for the current policy (Figure 2).

155 Thus, iterative SFT addresses the memory problem but leaves a core tension. Broad buffer coverage
 156 is needed for high dataset-level performance, but obtaining that coverage through fresh sampling
 157 can be expensive. Re-using data from old policy samples or other sources is the reasonable choice
 158 reduces this cost, but, in the worst cases it can lead to reward degradation. In the next section, we
 159 will aim to address this tension.

160 4 From Target Distribution to Algorithm

161 The preceding discussion suggests a simple design principle: we should retain a buffer of past
 162 successes, but avoid updating toward the buffer alone. Instead, the target distribution should involve
 163 the buffer where it provides high-reward signal and preserve the current policy elsewhere. We
 164 therefore propose to combine buffered successes with online samples from the current policy, then
 165 update toward this target through a new update that yields a one-step improvement guarantee.

166 Our algorithm, called Good Experience Maximization (GEM) is conceptually simple and consists
 167 of two main steps, detailed in Algorithm 1. First, it constructs P_M from the current buffer and the
 168 current policy (*Generate Good Experiences*). Second, rather than taking an unconstrained SFT step
 169 toward π_M , it follows a distributional flow from π_t to π_M (*Maximize Them*). This flow preserves the
 170 support benefits of the mixture target while allowing us to prove one-step improvement. We discuss
 171 each of these in turn.

172 4.1 Generate Good Experiences

173 We introduce a mixture distribution with weights $0 < \alpha < 1$,

$$P_M(\mathbf{x}, \mathbf{y}) = \alpha P_{\mathcal{B}}(\mathbf{x}, \mathbf{y}) + (1 - \alpha) \mu(\mathbf{x}) \pi_t(\mathbf{y}|\mathbf{x}), \quad (6)$$

174 which corresponds to line 5 in Algorithm 1. Intuitively, the first term concentrates mass on (\mathbf{x}, \mathbf{y})
 175 pairs that are high reward; the second term ensures that π_t covers \mathbf{x} outside of the support of \mathcal{B} .

Algorithm 1 GEM

- 1: Base LLM π_1 , buffer $\mathcal{B} \leftarrow \{\emptyset\}$
- 2: **for** policy iteration $t = 1, \dots, T$ **do**
- 3: Initialize training policy $\pi_\theta \leftarrow \pi_t$
- 4: **while** not converged **do**
- 5: Sample from the mixture $P_M(\mathbf{x}, \mathbf{y}) = \alpha P_{\mathcal{B}}(\mathbf{x}, \mathbf{y}) + (1 - \alpha)\mu(\mathbf{x})\pi_t(\mathbf{y}|\mathbf{x})$
- 6: Minimize divergence to mixture, $D_{\text{KL}}(P_M||\pi_\theta)$:

$$\hat{l}(\theta) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P_M} \left[-\log \pi_\theta(\mathbf{y}|\mathbf{x}) \right] \quad (5)$$

- 7: Annotate new policy samples with verifier $r_i = R(\mathbf{x}_i, \mathbf{y}_i)$
 - 8: Add good samples to buffer: $\mathcal{B} \leftarrow \mathcal{B} \cup \{(\mathbf{x}_i, \mathbf{y}_i, r_i, t)\}_{i=1}^n$
 - 9: **end while**
 - 10: $\pi_{t+1} \leftarrow \pi_\theta$
 - 11: **end for**
-

176 Define the conditional mixture policy as $\pi_M(\mathbf{y} | \mathbf{x}) = P_M(\mathbf{x}, \mathbf{y})\mu_M(\mathbf{x})$, where μ_M is the marginal
 177 of the mixture. Let $S_{\mathcal{B}}\{\mathbf{x} \in \mathcal{X} : \mu_{\mathcal{B}} > 0\}$ denote the support of the buffer’s marginal distribution.

178 **The mixture target reward is pointwise non-decreasing.** We will first show that this construction
 179 guarantees that expected reward does not decrease for any individual \mathbf{x} .

180 **Lemma 4.1.** *Suppose that, $\forall \mathbf{x} \in S_{\mathcal{B}}$, $\pi_{\mathcal{B}}$ outperforms π_t , such that: $\mathbb{E}_{\mathbf{y} \sim \pi_{\mathcal{B}}}[R(\mathbf{x}, \mathbf{y})] >$
 181 $\mathbb{E}_{\mathbf{y} \sim \pi_t}[R(\mathbf{x}, \mathbf{y})]$. Then, for π_M , for $\mathbf{x} \in S_{\mathcal{B}}$, $\mathbb{E}_{\mathbf{y} \sim \pi_{\mathcal{B}}}[R(\mathbf{x}, \mathbf{y})] > \mathbb{E}_{\mathbf{y} \sim \pi_t}[R(\mathbf{x}, \mathbf{y})]$. And, for $\mathbf{x} \notin S_{\mathcal{B}}$,
 182 $\mathbb{E}_{\mathbf{y} \sim \pi_{\mathcal{B}}}[R(\mathbf{x}, \mathbf{y})] = \mathbb{E}_{\mathbf{y} \sim \pi_t}[R(\mathbf{x}, \mathbf{y})]$. Therefore, $\mathbb{E}_{\pi_M}[R(\mathbf{x}, \mathbf{y})] \geq \mathbb{E}_{\pi_t}[(\mathbf{x}, \mathbf{y})]$.*

183 *Proof sketch.* For any \mathbf{x} , π_M is a convex combination of $\pi_{\mathcal{B}}$ and π_t , weighted by the marginals. For
 184 $\mathbf{x} \in S_{\mathcal{B}}$ the reward is the weighted average of $\mathbb{E}_{\mathbf{y} \sim \pi_{\mathcal{B}}}[R(\mathbf{x}, \mathbf{y})]$ and $\mathbb{E}_{\mathbf{y} \sim \pi_t}[R(\mathbf{x}, \mathbf{y})]$. The former is
 185 strictly greater by construction, so the average exceeds the latter. For $\mathbf{x} \notin S_{\mathcal{B}}$, π_M reduces to π_t . See
 186 Appendix A for details. \square

187 **Pointwise non-decreasing reward implies global improvement.** We now can consider the total
 188 expected reward across the *entire state distribution* $\mu(\mathbf{x})$. If a training algorithm perfectly matches the
 189 target conditional distribution, the new policy will converge to $\pi_M(\mathbf{y}|\mathbf{x})$. By integrating the pointwise
 190 improvement over $\mu(\mathbf{x})$, we show that this converged policy strictly improves global expected reward:
 191

192 **Lemma 4.2.** *Let $P_{\mathcal{B}} = \mu_{\mathcal{B}}(\mathbf{x})\pi_{\mathcal{B}}(\mathbf{y} | \mathbf{x})$. Assume $P_{\mathbf{x} \sim \mu}(\mathbf{x} \in S_{\mathcal{B}}) > 0$. Given the pointwise
 193 conditions from Lemma 4.1, the converged mixture policy π_M strictly outperforms the current policy
 194 π_t : Then, for any $\alpha \in (0, 1)$,*

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P_M} [R(\mathbf{x}, \mathbf{y})] > \mathbb{E}_{\mathbf{x} \sim \mu, \mathbf{y} \sim \pi_{\theta_t}} [R(\mathbf{x}, \mathbf{y})]. \quad (7)$$

195

196 *Proof sketch.* The expected reward under P_M can be decomposed by linearity as $\mathcal{J}(P_M) =$
 197 $\alpha \mathbb{E}_{P_{\mathcal{B}}}[R] + (1 - \alpha)\mathcal{J}(\pi_{\theta_t})$. By the assumption that \mathcal{B} outperforms the π_{θ_t} on its own support,
 198 $\mathbb{E}_{P_{\mathcal{B}}}[R] > \mathbb{E}_{\mathbf{x} \sim \mu, \mathbf{y} \sim \pi_{\theta_t}}[R]$. Since $\alpha > 0$, we get a strict inequality. See Appendix A for details. \square

199 4.2 Maximize Them

200 Theorem 4.2 guarantees that P_M has higher reward than π_t . To translate this statement into im-
 201 provement of π_{t+1} itself, we need π_{t+1} to be close to P_M , which holds only at convergence. With
 202 finite gradient steps, π_{t+1} may not yet match P_M , and individual steps along the way can decrease
 203 expected reward (Figure 1). In other words, Algorithm 1 guarantees per-iteration improvement only
 204 when *each inner loop runs to convergence*; per-update improvement is not guaranteed. We desire an
 205 update that preserves step-wise improvement toward P_M .

Algorithm 2 FlowUpdate($\mathcal{X}, \mathcal{Y}, \pi_\theta, m$)

- 1: Compute probabilities $p = \pi_\theta(\cdot | \mathcal{X})$
 - 2: Compute probability velocity $\dot{p} = m - p$
 - 3: Map to parameter space $L = \dot{p} \oslash p$
 - 4: **return** L
-

206 **Gradient descent can leave the mixture path.** Observe that we can specify a path in distribution
207 space from the current policy to the mixture target. For convenience of notation, we write the current
208 policy $p = \pi_\theta$ and the target distribution supported on correct answers as b . Rewriting Equation (6),
209 note that this defines a convex combination of the two distributions in probability space:

$$m(\alpha) = (1 - \alpha)p + \alpha b, \alpha \in [0, 1].$$

210 Moving along this path corresponds to a straight line in the probability simplex (an m -geodesic [Amari,
211 2016]). However, although b represents the ideal destination on the probability simplex, standard
212 gradient descent in parameter space does not necessarily recover this path. Instead, standard gradi-
213 ent steps inherently follow trajectories dictated by the parameter geometry, which can violate the
214 guarantee of monotonic performance improvement [Kakade, 2001] (which we observe in Figures 1
215 and 2).

216 **A probability-space flow recovers the path.** To restore this guarantee, we explicitly force the
217 policy to move along the straight line path by constructing a continuous-time flow (Algorithm 2).
218 Note that the ideal velocity in probability space along this path is simply the vector pointing from
219 the current policy to the target distribution: $\dot{p} = \alpha - p$. To realize the velocity \dot{p} in probability
220 space through continuous steps in parameter space, we use the chain rule to derive the corresponding
221 parameter velocity $\dot{\theta}$. Because the expected reward is a linear function of the output probabilities, it
222 strictly interpolates between the performance of the current policy and the target policy when moving
223 in a straight line toward a more performant target distribution.

224 **With FlowUpdate, GEM theoretically improves performance at every step.** The advantage of
225 this framing is it guarantees performance gains at *every step*, which we will now show.

226 **Proposition 4.3.** *Let $\mathcal{J}(\theta)$ be the expected reward of a policy parameterized by its logits θ . Assume*
227 *$\dot{\theta}$ is defined component-wise as $\dot{\theta}_i := \frac{\dot{p}_i}{p_i} - \frac{\dot{p}_v}{p_v}$, which realizes $\dot{p} = m - p$. If m has non-negative*
228 *expected advantage, $\sum_j m_j A_j \geq 0$, then: $\frac{d\mathcal{J}}{dt} \geq 0$.*

229 *Proof sketch.* We compute the time derivative of performance via the chain rule: $\frac{d\mathcal{J}}{dt} = \sum_{i=1}^{v-1} \frac{\partial \mathcal{J}}{\partial \theta_i} \dot{\theta}_i$.
230 Substituting the policy gradient $\frac{\partial \mathcal{J}}{\partial \theta_i} = p_i A_i$ and our explicit update, we get $\sum_{i=1}^{v-1} p_i A_i \left(\frac{\dot{p}_i}{p_i} - \frac{\dot{p}_v}{p_v} \right)$,
231 which simplifies to $\sum_{i=1}^v A_i \dot{p}_i$. Plugging in the trajectory $\dot{p}_i = m_i - p_i$, we get $\frac{d\mathcal{J}}{dt} = \sum_i (m_i - p_i) A_i$.
232 Since the expected advantage of the current policy is always 0, the derivative exactly equals $\sum_i m_i A_i$,
233 which is non-negative. Therefore, $\frac{d\mathcal{J}}{dt} \geq 0$. We provide the full derivation in Section A. \square

234 4.3 Practical Desiderata

235 Given the generality of GEM, we now discuss a few key design choices: what to retain, how to
236 sample, and how to amortize online generation.

237 **What to retain.** In the binary-reward setting, retention is straightforward. Storing correct genera-
238 tions automatically satisfies the condition in Theorem 4.2 because $\mathbb{E}_{\mathcal{B}}[R] = 1$, whereas $E_{\pi_i}[R] \leq 1$.
239 For a fixed capacity buffer, a natural eviction rule is FIFO, which discards the most stale data first.
240 However, more complex eviction rules may prioritize maintaining diversity or difficulty over \mathcal{X} , the
241 latter of which would likely better maintain the support condition (especially if the buffer capacity is
242 small relative to the number of questions in the training dataset).

243 **How to sample.** $P_{\mathcal{B}}$ does not require a specific sampling scheme beyond coverage, but different
244 choices can have benefit in practice. The simplest approach is to train on uniformly sampled entries;
245 however, uniform sampling introduces two sources of bias. First, samples that enter the buffer earlier
246 are disproportionately sampled. Thus, a simple fix is to sample with a recency bias according to an
247 exponentially decaying weight. Second, easy problems (high $\pi_\theta(\mathcal{A}^*(\mathbf{x}) | \mathbf{x})$) contribute many more
248 correct entries to the buffer than hard problems. As a result, training is dominated by examples the

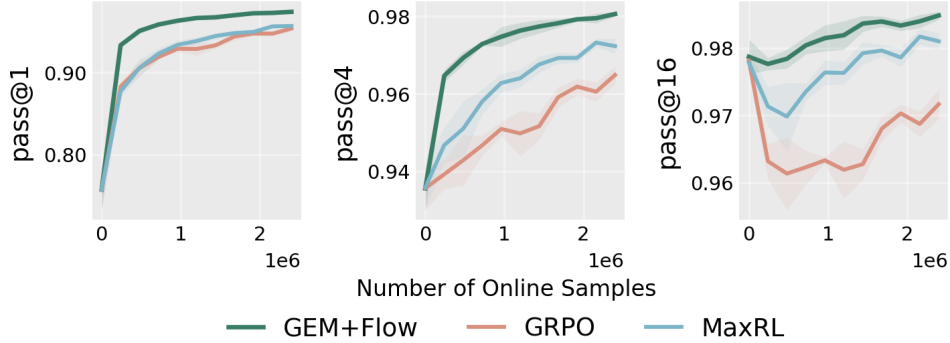


Figure 3: Number of online samples (generations) vs. pass@k for GEM+Flow, GRPO [Shao et al., 2024], and MaxRL [Tajwar et al., 2026] on the MNIST bandit task. Online samples denote newly generated model outputs used for training, excluding reused buffer samples. Curves show mean performance across 3 seeds, with shaded regions indicating standard deviation. GEM+Flow reaches higher pass@k with fewer online generations than GRPO or MaxRL, demonstrating improved sample efficiency.

249 model has already solved, providing little learning signal. We address this issue with *inverse-pass-rate*
 250 $(1/p)$ sampling. For each problem, define the empirical pass rate from the buffer as $\hat{p}_x = \frac{c_x}{n_x}$, where
 251 n_x is the total number of entries for x and c_x is the number of correct ones. The sampling weight
 252 is $w(x) = \frac{1}{\hat{p}_x}$. Under these weights, the expected number of training samples from each prompt is
 253 equalized regardless of difficulty. ¹

254 **How to amortize online generation.** Each iteration
 255 of Algorithm 1 requires some number of online rollouts
 256 from π_θ (line 5). Because generation is costly, amortizing
 257 it across multiple steps is desirable. One approach is to
 258 occasionally skip generating online samples and update
 259 purely using a static buffer. This version holds given the
 260 buffer-reward condition is maintained, which occurs when
 261 the new policy is updated with a small learning rate or is
 262 updated using a trust region. Figure 3 shows how GEM
 263 can make use of additional updates on a ViT on MNIST
 264 bandit task. Empirically, Figure 3 shows that GEM im-
 265 proves pass@1 and pass@4 monotonically while avoid-
 266 ing the larger early pass@16 degradation of GRPO and
 267 MaxRL. Moreover, GEM reaches a given pass@k level
 268 with substantially fewer newly-collected samples, indicat-
 269 ing improved sample efficiency.

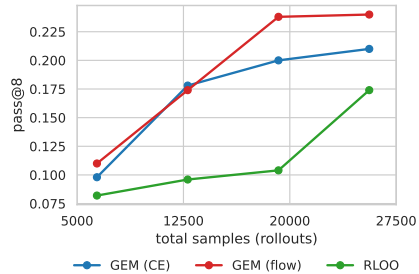


Figure 4: Eval Pass@8 of Llama on MATH. GEM achieves fast improvements, with the flow update providing slightly better learning

270 5 A Practical Algorithm

271 We introduced the flow update due to the phenomenon that standard SFT can curve away from the
 272 straight-line mixture path in probability space. However, the current policy may actually lie in a
 273 benign region of the simplex, such as in the center of those shown in Figure 2. In these cases, the
 274 SFT update may approximate the flow update well, moving more directly to the target distribution
 275 without reward degradation. It is difficult to know a priori which regime will hold during LM
 276 post-training, but, if the mixture target can often be optimized with ordinary SFT, then GEM can
 277 inherit the practical advantages of existing supervised-training pipelines. We therefore test a practical
 278 SFT version of GEM that uses the same online replay construction but replaces the flow update with
 279 standard cross-entropy training. Critically, this variant does not train on a fixed dataset: it repeatedly
 280 generates new completions and trains on the evolving replay distribution.

¹This in fact recovers recent works [Tajwar et al., 2026] in the setting where the buffer only maintains storage of the most recent policy’s generations.

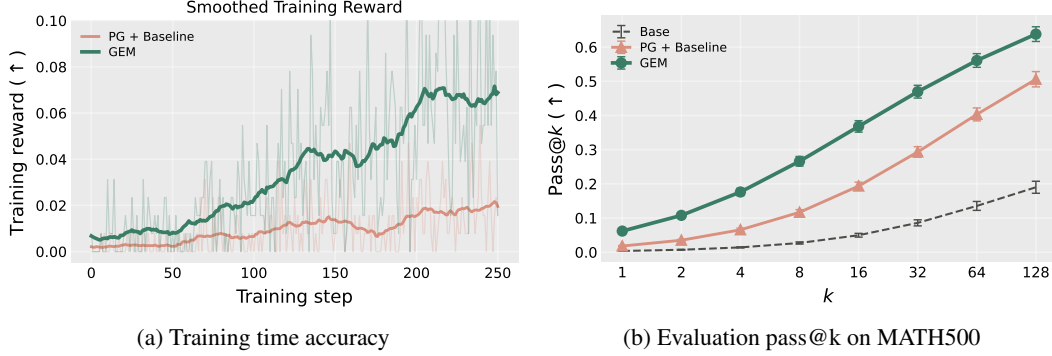


Figure 5: *GEM* matches or excels performance of policy gradient (PG with mean reward baseline) and base model (Llama 3.2 3B Base), particularly excelling when rewards are difficult to find, demonstrating its ability to make use of previously seen high reward generations (on MATH, while the base model has nearly zero pass@1 rates).

281 A central motivation for *GEM* is that policy-gradient training can waste rare successes: once a
 282 high-reward output is sampled, the update may not preserve the specific solution for future learning.
 283 In the dense-reward setting, many outputs provide a useful signal for learning; in the sparse-reward
 284 setting, successful outputs are rare. We therefore train on MATH [Hendrycks et al., 2021], and
 285 evaluate using MATH500 [Lightman et al., 2023]. We first validate that both versions of *GEM*
 286 can indeed provide performance and efficiency gains on LM post-training. We then investigate the
 287 practical version of *GEM*. For these experiments, we use Llama-3.2-3B Base [Grattafiori et al., 2024].
 288 For additional experimental details, see Appendix B.3.

289 **Both variants can provide reward improvement on LMs.** We first demonstrate that *GEM* with
 290 and without the flow update can outperform a standard policy gradient post-training method (RLOO)
 291 using Figure 4 demonstrates that *GEM* rapidly improves, with the flow update providing marginally
 292 better learning.

293 **More broadly, *GEM* without the flow update outperforms RL.** We now turn to the practical
 294 algorithm and investigate its performance. Figure 5 shows that our algorithm outperforms the
 295 standard policy gradient in sparse reward settings. These results support that maintaining an explicit
 296 memory of past successes is particularly beneficial in such “hard discovery” settings.

297 6 A Unifying View: Policy Optimization is Iterative SFT

298 In this section we expand on the note that all existing work around policy optimization for large
 299 generative models, which operate in the contextual bandit setting, is more simply viewed through the
 300 lens of iterative minimization of a forward KL divergence (i.e. “SFT update”) towards *some* target
 301 distribution P ,

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P} [- \nabla_{\theta} \log \pi_{\theta}(\mathbf{y} | \mathbf{x})] .$$

302 To see this, we note that for any positive reward the policy gradient can be written as,

$$\begin{aligned} \nabla_{\theta} \mathcal{J}_{\text{PG}}(\pi_{\theta}) &= \mathbb{E}_{\mu(\mathbf{x})} \left[\mathbb{E}_{\pi_{\theta}(\mathbf{y} | \mathbf{x})} [R(\mathbf{x}, \mathbf{y}) \nabla_{\theta} \log \pi_{\theta}(\mathbf{y} | \mathbf{x})] \right] \\ &= \mathbb{E}_{\mu(\mathbf{x})} \left[\int \pi_{\theta}(\mathbf{y} | \mathbf{x}) R(\mathbf{x}, \mathbf{y}) \nabla_{\theta} \log \pi_{\theta}(\mathbf{y} | \mathbf{x}) \right], \\ &= \mathbb{E}_{\mu(\mathbf{x})} \left[v(\mathbf{x}, \pi_{\theta}) \int \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x}) R(\mathbf{x}, \mathbf{y})}{v(\mathbf{x}, \pi_{\theta})} \nabla_{\theta} \log \pi_{\theta}(\mathbf{y} | \mathbf{x}) \right], \\ &= \mathbb{E}_{\mu(\mathbf{x})} \left[v(\mathbf{x}, \pi_{\theta}) \mathbb{E}_{q_{\theta}(\mathbf{y} | \mathbf{x})} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{y} | \mathbf{x})] \right], \end{aligned}$$

303 where $v(\mathbf{x}, \pi_{\theta}) = \mathbb{E}_{\pi_{\theta}(\mathbf{y} | \mathbf{x})} [R(\mathbf{x}, \mathbf{y})]$ is the expected reward of π_{θ} given input \mathbf{x} . The above is
 304 equivalent to a per-question SFT update toward the target,

$$q_{\theta}(\mathbf{y} | \mathbf{x}) = \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x}) R(\mathbf{x}, \mathbf{y})}{v(\mathbf{x}, \pi_{\theta})} .$$

Method	Target Distribution	Per Question Scaling	Sampling Req	Update Steps
Policy Gradient	$\propto \pi_t(\mathbf{y} \mathbf{x})R(\mathbf{x}, \mathbf{y})$	1	Online	One/Few
EM	$\propto \pi_t(\mathbf{y} \mathbf{x})R(\mathbf{x}, \mathbf{y})$	$1/v(\mathbf{x}, \pi_t)$	Online	One/Few
Filtered SFT	$\propto \pi_t(\mathbf{y} \mathbf{x})R(\mathbf{x}, \mathbf{y})$	1	High	Many
GEM	$\propto \pi_{t,t-1,\dots}(\mathbf{y} \mathbf{x})R(\mathbf{x}, \mathbf{y})$	–	Online	Many

Table 1: Policy optimization methods at policy iteration t

305 The above describes any set of policy gradient updates [Williams, 1992, Schulman et al., 2017, Guo
306 et al., 2025]. It is interesting to note that the per-question gradient is scaled by a factor proportion to
307 the performance on that question, $v(\mathbf{x}, \pi_\theta)$. Correspondingly, a set of methods based on expectation
308 maximization to directly increase the likelihood of success [Dayan and Hinton, 1997, Abdolmaleki
309 et al., 2018b, 2024, Tajwar et al., 2026] do weighted SFT updates which effectively has the effect of
310 canceling out this scaling factor,

$$\begin{aligned} \nabla \mathcal{J}_{\text{EM}}(\pi_\theta) &= \mathbb{E}_{\mu(\mathbf{x})} \left[\frac{1}{v(\mathbf{x}, \pi_\theta)} \mathbb{E}_{\pi_\theta(\mathbf{y}|\mathbf{x})} [R(\mathbf{x}, \mathbf{y}) \nabla_\theta \log \pi_\theta(\mathbf{y}|\mathbf{x})] \right] \\ &= \mathbb{E}_{\mu(\mathbf{x})} \left[\mathbb{E}_{q_\theta(\mathbf{y}|\mathbf{x})} [\nabla_\theta \log \pi_\theta(\mathbf{y}|\mathbf{x})] \right]. \end{aligned}$$

311 It is worth noting that q_θ is precisely also the distribution we would get in the binary reward case from
312 filtering. That is, iterative filtered SFT, REINFORCE, and EM-based RL methods are optimizing the
313 same expected gradient (up to a per-question scaling factor). We summarize this in Table 1, along
314 with our algorithm GEM.

315 **Recovery of Existing Methods from GEM.** We can recover various existing methods from GEM
316 by picking a specific choice of buffer sampling distribution and update steps. For instance, to recover
317 **MaxRL** [Tajwar et al., 2026] or EM [Dayan, 1993] requires setting buffer capacity to $C = BK$
318 (buffer holds exactly one batch, so data is never reused across collection rounds), $G = 1$, mini-batches
319 are drawn *without replacement* so that every correct rollout in the buffer is used exactly once, and
320 the sampling weights are $1/\hat{p}_q$. Under these settings the effective contribution of each prompt \mathbf{x} to
321 the loss is the sum over all its correct rollouts, each weighted by $1/\hat{p}_q$. Because $\hat{p}_q = c_q/n_q$ and
322 there are c_q correct entries, the total weight for prompt \mathbf{x} is $c_q \cdot (n_q/c_q) = n_q$, which is constant
323 across prompts. This replicates the MaxRL estimator: every correct rollout contributes equally to
324 the gradient regardless of the prompt’s difficulty, matching the uniform-over-prompts weighting
325 implicit in the MaxRL gradient. Our method extends MaxRL by (a) maintaining a persistent buffer
326 ($C \gg BK$) so correct traces are reused across collection rounds, and (b) performing $G > 1$ gradient
327 steps per collection phase.

328 7 Conclusion

329 We introduced Good Experience Maximization (GEM), a post-training framework that retains and
330 reuses successful generations while continuing to collect new online samples. Our target-matching
331 view shows key tensions between policy gradient and iterative SFT. GEM addresses this tension by
332 training toward a mixture of historical and newly generated successes. We proved that the GEM
333 target improves expected reward at convergence and derived a flow-based update that guarantees
334 one-step policy improvement by moving along a straight-line path in probability space. Empirically,
335 GEM improves sample efficiency and performs particularly well in sparse-reward settings. Future
336 work should scale these experiments further, study practical buffer design, and characterize when
337 standard SFT is sufficient versus when the flow update is needed. Additionally, we note that we
338 do not make use of *negative* samples [Abdolmaleki et al., 2024] in our algorithm. These are useful
339 because they can help guide search and can provide rich information depending on the structure (e.g.,
340 a negative sample of two options is equally as informative as a positive one).

341 References

- 342 Abbas Abdolmaleki, Jost Tobias Springenberg, Jonas Degraeve, Steven Bohez, Yuval Tassa, Dan
343 Belov, Nicolas Heess, and Martin Riedmiller. Relative entropy regularized policy iteration. *arXiv*
344 *preprint arXiv:1812.02256*, 2018a.
- 345 Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and Martin
346 Riedmiller. Maximum a posteriori policy optimisation. *arXiv preprint arXiv:1806.06920*, 2018b.
- 347 Abbas Abdolmaleki, Bilal Piot, Bobak Shahriari, Jost Tobias Springenberg, Tim Hertweck, Rishabh
348 Joshi, Junhyuk Oh, Michael Bloesch, Thomas Lampe, Nicolas Heess, et al. Learning from negative
349 feedback, or positive feedback or both. *arXiv preprint arXiv:2410.04166*, 2024.
- 350 Shun-ichi Amari. *Information geometry and its applications*. Springer, 2016.
- 351 Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob
352 McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay.
353 *Advances in neural information processing systems*, 30, 2017.
- 354 David Brandfonbrener, William F Whitney, Rajesh Ranganath, and Joan Bruna. Quantile filtered
355 imitation learning. *arXiv preprint arXiv:2112.00950*, 2021.
- 356 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
357 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
358 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 359 Peter Dayan. Improving generalization for temporal difference learning: The successor representation.
360 *Neural Computation*, 1993.
- 361 Peter Dayan and Geoffrey E Hinton. Using expectation-maximization for reinforcement learning.
362 *Neural Computation*, 9(2):271–278, 1997.
- 363 Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao,
364 Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative
365 foundation model alignment. *Transactions on Machine Learning Research*, 2023.
- 366 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
367 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of
368 models. *arXiv preprint arXiv:2407.21783*, 2024.
- 369 Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek
370 Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, et al. Reinforced self-training
371 (ReST) for language modeling. *arXiv preprint arXiv:2308.08998*, 2023.
- 372 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu
373 Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through reinforce-
374 ment learning. *Nature*, 645(8081):633–638, 2025.
- 375 Alex Havrilla, Yuqing Du, Sharath Chandra Raparthi, Christoforos Nalmpantis, Jane Dwivedi-Yu,
376 Maksym Zhuravinskyi, Eric Hambro, Sainbayar Sukhbaatar, and Roberta Raileanu. Teaching large
377 language models to reason with reinforcement learning. *arXiv preprint arXiv:2403.04642*, 2024.
- 378 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
379 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
380 *preprint arXiv:2103.03874*, 2021.
- 381 Sham M Kakade. A natural policy gradient. *Advances in Neural Information Processing Systems*, 14,
382 2001.
- 383 Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman,
384 Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in
385 open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.

- 386 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
387 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
388 *arXiv:2305.20050*, 2023.
- 389 Long-Ji Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching.
390 *Machine Learning*, 8(3–4):293–321, 1992.
- 391 Marcelo G Mattar and Nathaniel D Daw. Prioritized memory access explains planning and hippocam-
392 pal replay. *Nature neuroscience*, 21(11):1609–1617, 2018.
- 393 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
394 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint*
395 *arXiv:1312.5602*, 2013.
- 396 Ida Momennejad, A Ross Otto, Nathaniel D Daw, and Kenneth A Norman. Offline replay supports
397 planning in human reinforcement learning. *elife*, 7:e32548, 2018.
- 398 Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured
399 prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on*
400 *artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings,
401 2011.
- 402 Anna C Schapiro, Elizabeth A McDevitt, Timothy T Rogers, Sara C Mednick, and Kenneth A
403 Norman. Human hippocampal replay during rest prioritizes weakly learned information and
404 predicts memory performance. *Nature communications*, 9(1):3920, 2018.
- 405 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
406 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 407 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
408 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathemat-
409 ical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 410 Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Xavier Garcia, Peter J
411 Liu, James Harrison, Jaehoon Lee, Kelvin Xu, et al. Beyond human data: Scaling self-training for
412 problem-solving with language models. *arXiv preprint arXiv:2312.06585*, 2023.
- 413 Yi Su, Dian Yu, Linfeng Song, Juntao Li, Haitao Mi, Zhaopeng Tu, Min Zhang, and Dong Yu.
414 Crossing the reward bridge: Expanding rl with verifiable rewards across diverse domains. *arXiv*
415 *preprint arXiv:2503.23829*, 2025.
- 416 Fahim Tajwar, Guanning Zeng, Yuer Zhou, Yuda Song, Daman Arora, Yiding Jiang, Jeff Schneider,
417 Ruslan Salakhutdinov, and Haiwen Feng. Maximum likelihood reinforcement learning. *arXiv*
418 *preprint arXiv:2602.02710*, 2026.
- 419 Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung
420 Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in
421 starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- 422 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
423 learning. *Machine learning*, 8(3):229–256, 1992.
- 424 Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi
425 Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on
426 imagenet. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages
427 558–567, 2021.
- 428 Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does
429 reinforcement learning really incentivize reasoning capacity in llms beyond the base model?
430 *Advances in Neural Information Processing Systems*, 2025.
- 431 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with
432 reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.

433 A Deferred Proofs

434 A.1 Proofs from Section 4

435 *Proof of Lemma 4.1.* By the definition provided in the sketch, the mixture policy π_M is a convex
 436 combination of π_B and π_t for any \mathbf{x} . We can write the probability of \mathbf{y} in \mathbf{x} under π_M as $\pi_M(\mathbf{y}|\mathbf{x}) =$
 437 $\alpha(\mathbf{x})\pi_B(\mathbf{y}|\mathbf{x}) + (1 - \alpha(\mathbf{x}))\pi_t(\mathbf{y}|\mathbf{x})$, where $\alpha(\mathbf{x})$. First consider the case where $\mathbf{x} \in S_B$. In this
 438 region, $\alpha(\mathbf{x}) > 0$. The expected reward at \mathbf{x} under π_M is given by:

$$\begin{aligned} \mathbb{E}_{\mathbf{y} \sim \pi_M} [R(\mathbf{x}, \mathbf{y})] &= \sum_{\mathbf{y} \in \mathcal{A}} \pi_M(\mathbf{y}|\mathbf{x}) R(\mathbf{x}, \mathbf{y}) \\ &= \sum_{\mathbf{y} \in \mathcal{A}} (\alpha(\mathbf{x})\pi_B(\mathbf{y}|\mathbf{x}) + (1 - \alpha(\mathbf{x}))\pi_t(\mathbf{y}|\mathbf{x})) R(\mathbf{x}, \mathbf{y}) \\ &= \alpha(\mathbf{x}) \mathbb{E}_{\mathbf{y} \sim \pi_B} [R(\mathbf{x}, \mathbf{y})] + (1 - \alpha(\mathbf{x})) \mathbb{E}_{\mathbf{y} \sim \pi_t} [R(\mathbf{x}, \mathbf{y})]. \end{aligned}$$

439 By the assumption that $\mathbb{E}_{\mathbf{y} \sim \pi_B} [R(\mathbf{x}, \mathbf{y})] > \mathbb{E}_{\mathbf{y} \sim \pi_t} [R(\mathbf{x}, \mathbf{y})]$ for all $\mathbf{x} \in S_B$, and since $\alpha(\mathbf{x}) > 0$,
 440 the weighted average must be strictly greater than the target expectation. Thus, $\mathbb{E}_{\mathbf{y} \sim \pi_M} [R(\mathbf{x}, \mathbf{y})] >$
 441 $\mathbb{E}_{\mathbf{y} \sim \pi_t} [R(\mathbf{x}, \mathbf{y})]$ for all $\mathbf{x} \in S_B$.

442 Next, consider the case where $\mathbf{x} \notin S_B$. For these, the mixture weight $\alpha(\mathbf{x})$ reduces to 0. Consequently,
 443 the mixture policy simplifies to $\pi_M(\mathbf{y}|\mathbf{x}) = \pi_t(\mathbf{y}|\mathbf{x})$. It follows immediately that the expectations
 444 are equal: $\mathbb{E}_{\mathbf{y} \sim \pi_M} [R(\mathbf{x}, \mathbf{y})] = \mathbb{E}_{\mathbf{y} \sim \pi_t} [R(\mathbf{x}, \mathbf{y})]$.

445 Finally, to arrive at the global inequality, we take the expectation over $d(\mathbf{x})$. The total expected
 446 reward is:

$$\begin{aligned} \mathbb{E}_{\pi_M} [R(\mathbf{x}, \mathbf{y})] &= \int_{S_B} d(\mathbf{x}) \mathbb{E}_{\mathbf{y} \sim \pi_M} [R(\mathbf{x}, \mathbf{y})] d\mathbf{x} + \int_{S \setminus S_B} d(\mathbf{x}) \mathbb{E}_{\mathbf{y} \sim \pi_M} [R(\mathbf{x}, \mathbf{y})] d\mathbf{x} \\ &\geq \int_{S_B} d(\mathbf{x}) \mathbb{E}_{\mathbf{y} \sim \pi_t} [R(\mathbf{x}, \mathbf{y})] d\mathbf{x} + \int_{S \setminus S_B} d(\mathbf{x}) \mathbb{E}_{\mathbf{y} \sim \pi_t} [R(\mathbf{x}, \mathbf{y})] d\mathbf{x} \\ &= \mathbb{E}_{\pi_t} [R(\mathbf{x}, \mathbf{y})]. \end{aligned}$$

447 The inequality holds because the integrand is strictly greater on S_B and equal elsewhere. Thus,
 448 $\mathbb{E}_{\pi_M} [R(\mathbf{x}, \mathbf{y})] \geq \mathbb{E}_{\pi_t} [R(\mathbf{x}, \mathbf{y})]$. \square

449 *Proof of Lemma 4.2.* We consider the total expected reward under the converged mixture distribution
 450 P_M . By definition, the mixture distribution is a convex combination of the buffer distribution and the
 451 current policy's distribution, expressed as $P_M = \alpha P_B + (1 - \alpha)P_{\pi_t}$, where $P_{\pi_t} = \mu(\mathbf{x})\pi_{\theta_t}(\mathbf{y}|\mathbf{x})$
 452 represents the distribution under the current policy. Let $\mathcal{J}(P) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P} [R(\mathbf{x}, \mathbf{y})]$ denote the
 453 expected reward under a given distribution P . We can expand the expectation for the mixture policy
 454 by substituting the definition of P_M :

$$\mathcal{J}(P_M) = \int_S \int_{\mathcal{A}} R(\mathbf{x}, \mathbf{y}) [\alpha P_B(\mathbf{x}, \mathbf{y}) + (1 - \alpha)P_{\pi_t}(\mathbf{x}, \mathbf{y})] d\mathbf{y} d\mathbf{x}.$$

455 By the linearity of the integral, we can decompose this expression into two distinct components
 456 corresponding to the reward contribution of the buffer and the current policy:

$$\mathcal{J}(P_M) = \alpha \iint_{S \times \mathcal{A}} R(\mathbf{x}, \mathbf{y}) P_B(\mathbf{x}, \mathbf{y}) d\mathbf{y} d\mathbf{x} + (1 - \alpha) \iint_{S \times \mathcal{A}} R(\mathbf{x}, \mathbf{y}) P_{\pi_t}(\mathbf{x}, \mathbf{y}) d\mathbf{y} d\mathbf{x}.$$

457 This expression simplifies to a weighted sum of the expected rewards:

$$\mathcal{J}(P_M) = \alpha \mathbb{E}_{P_B} [R] + (1 - \alpha) \mathcal{J}(P_{\pi_t}).$$

458 Based on the pointwise conditions from Lemma 4.1, we assume that the expected reward of the
 459 buffer distribution strictly exceeds that of the current policy over the shared support, such that
 460 $\mathbb{E}_{P_B} [R] > \mathcal{J}(P_{\pi_t})$. Given that the mixing weight α is strictly greater than zero, we can apply this
 461 strict inequality to our decomposition:

$$\mathcal{J}(P_M) > \alpha \mathcal{J}(P_{\pi_t}) + (1 - \alpha) \mathcal{J}(P_{\pi_t}).$$

462 Factoring out the term $\mathcal{J}(P_{\pi_t})$ on the right-hand side, we obtain:

$$\mathcal{J}(P_M) > (\alpha + 1 - \alpha)\mathcal{J}(P_{\pi_t}) = \mathcal{J}(P_{\pi_t}).$$

463 Thus, we have shown that for any mixture weight $\alpha \in (0, 1)$, the expected reward under the converged
464 mixture policy P_M is strictly greater than the expected reward under the current policy π_{θ_t} . \square

465 *Proof of Proposition 4.3.* To map the probability simplex to the parameter space θ , we select a
466 reference action K . The parameters θ_i for $i \in \{1, \dots, K-1\}$ are defined by the log-ratio:

$$\theta_i = \ln(p_i) - \ln(p_K)$$

467 subject to the simplex constraint $p_K = 1 - \sum_{j=1}^{K-1} p_j$. We first derive the Jacobian $J_{ij} = \frac{\partial \theta_i}{\partial p_j}$ to
468 understand how changes in probability space affect the parameters. For the diagonal case where
469 $i = j$:

$$\frac{\partial \theta_i}{\partial p_i} = \frac{\partial}{\partial p_i}(\ln p_i - \ln p_K) = \frac{1}{p_i} - \frac{1}{p_K} \frac{\partial p_K}{\partial p_i} = \frac{1}{p_i} + \frac{1}{p_K}$$

470 For the off-diagonal case where $i \neq j$:

$$\frac{\partial \theta_i}{\partial p_j} = \frac{\partial}{\partial p_j}(\ln p_i - \ln p_K) = 0 - \frac{1}{p_K} \frac{\partial p_K}{\partial p_j} = \frac{1}{p_K}$$

471 Thus, the general form of the Jacobian is $\frac{\partial \theta_i}{\partial p_j} = \frac{\mathbb{1}_{\{i=j\}}}{p_i} + \frac{1}{p_K}$. Using the chain rule, the velocity of
472 the parameters $\dot{\theta}_i$ given a probability velocity \dot{p} is:

$$\dot{\theta}_i = \sum_{j=1}^{K-1} \frac{\partial \theta_i}{\partial p_j} \dot{p}_j = \frac{\dot{p}_i}{p_i} + \frac{1}{p_K} \sum_{j=1}^{K-1} \dot{p}_j.$$

473 Since the probabilities must sum to 1, the sum of their derivatives is zero, implying $\sum_{j=1}^{K-1} \dot{p}_j = -\dot{p}_K$.
474 Substituting this into the equation above yields the explicit parameter update:

$$\dot{\theta}_i = \frac{\dot{p}_i}{p_i} - \frac{\dot{p}_K}{p_K}$$

475 Write the time derivative of the expected reward \mathcal{J} using the chain rule over the policy parameters θ
476 as:

$$\frac{d\mathcal{J}}{dt} = \sum_{i=1}^{K-1} \frac{\partial \mathcal{J}}{\partial \theta_i} \dot{\theta}_i.$$

477 Recall the standard result for the gradient of the expected reward with respect to softmax logits:
478 $\frac{\partial \mathcal{J}}{\partial \theta_i} = p_i(Q_i - \mathbb{E}_{a \sim p}[Q(s, a)]) = p_i A_i$, where A_i is the advantage of action i . Substituting this
479 gradient and the definition of the logit dynamics $\dot{\theta}_i$ into the chain rule expression, we obtain:

$$\frac{d\mathcal{J}}{dt} = \sum_{i=1}^{K-1} p_i A_i \left(\frac{\dot{p}_i}{p_i} - \frac{\dot{p}_K}{p_K} \right).$$

480 Distributing the terms inside the summation yields:

$$\frac{d\mathcal{J}}{dt} = \sum_{i=1}^{K-1} \dot{p}_i A_i - \frac{\dot{p}_K}{p_K} \sum_{i=1}^{K-1} p_i A_i.$$

481 By the definition of the advantage function, the expected advantage under the current policy is zero,
482 i.e., $\sum_{i=1}^K p_i A_i = 0$. This implies that $\sum_{i=1}^{K-1} p_i A_i = -p_K A_K$. Substituting this identity into the
483 second term of our derivative expression, we have:

$$\frac{d\mathcal{J}}{dt} = \sum_{i=1}^{K-1} \dot{p}_i A_i - \frac{\dot{p}_K}{p_K} (-p_K A_K) = \sum_{i=1}^{K-1} \dot{p}_i A_i + \dot{p}_K A_K = \sum_{i=1}^K \dot{p}_i A_i.$$

484 We now plug in the trajectory of the probability distribution, $\dot{p}_i = m_i - p_i$. Substituting this into the
485 simplified summation gives:

$$\frac{d\mathcal{J}}{dt} = \sum_{i=1}^K (m_i - p_i) A_i = \sum_{i=1}^K m_i A_i - \sum_{i=1}^K p_i A_i.$$

486 As established previously, $\sum_{i=1}^K p_i A_i = 0$. Therefore, the derivative of the performance objective
487 reduces exactly to the expected advantage of the target distribution m under the current policy:

$$\frac{d\mathcal{J}}{dt} = \sum_{i=1}^K m_i A_i.$$

488 By the assumption that $\sum_j m_j A_j \geq 0$, it follows immediately that $\frac{d\mathcal{J}}{dt} \geq 0$.

489

□

Hyperparameter	Default Value
Learning rate	10^{-3}
Optimizer	Adam
Batch size	32
Rollouts per context	4
Total sample budget	50,000

Table 2: Default training hyperparameters for all MNIST experiments.

Parameter	Value
Base model	Llama-3.2-3B
LoRA rank	32
Train dataset	MATH
Eval dataset	MATH-500
Reward	binary
Prompts per batch	16
Generations per prompt	8
Max response length	2048
Training steps	400
Optimizer	Adam ($\beta_1=0.9, \beta_2=0.95$)
Grad updates per data collection step	1
Learning rate	1×10^{-5}
Rollout temp / top- p	1.0 / 1.0
Validation temp / top- p	1.0 / 0.95
Compute backend	Tinker (managed)

Table 3: Setup and data for math reasoning.

490 B Experiment Details

491 B.1 Pass@k Calculation

492 We use bootstrapping low variance unbiased estimator for estimating pass@k from samples [Chen
493 et al., 2021]. This estimator is calculated by generating $n \geq k$ samples per task and counting the
494 number of correct samples $c(\mathbf{x})$ among the n samples. The pass@k is therefore:

$$\text{pass@k} := \mathbb{E}_{\mathbf{x} \sim \mu} \left[1 - \frac{\binom{n-c(\mathbf{x})}{k}}{\binom{n}{k}} \right].$$

495 B.2 MNIST

496 We frame MNIST as a contextual bandit task. Each problem \mathbf{x} represents the context. It is a
497 28×28 grayscale image (flattened to \mathbb{R}^{784} , normalized to $[0, 1]$). Each answer in the answer space
498 $\mathcal{A} = \{0, \dots, C - 1\}$ (with $C = 10$ by default) represents an arm. The reward is binary, such that
499 $r(\mathbf{y}, y) = \mathbf{1}[y = y]$, where y is the true digit label. The MNIST training set contains 60,000 images
500 and the test set 10,000 images.

501 **Hyperparameters.** The results in Figure 3 were generated with the hyperparameters in Table 2.
502 All algorithms were trained using ViT architecture [Yuan et al., 2021] with 6 layers and 4 heads with
503 patch size 4. GEM used $\alpha = 0.5$, buffer capacity of 4096, $1/p$ sampling, 5 inner updates per outer
504 loop.

505 **Compute Resources.** Each run takes around 4-5 hours and requires a single L40.

506 B.3 Math Reasoning

507 **Hyperparameters.** Table ?? lists the hyperparameters for the mathematicla reasoning tasks.

508 **NeurIPS Paper Checklist**

509 **1. Claims**

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

512 Answer: [Yes]

513 Justification: The abstract and/or introduction clearly states the claims made.

514 **2. Limitations**

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

516 Answer: [Yes]

517 Justification: Throughout the paper, we explicitly mention failure cases and limitations of the
518 methods developed. We mention assumptions where possible and address the scope of the
519 work. In particular, because we are compute-limited, we could only run a few larger-scale
520 experiments. However, we provide supporting evidence in simpler settings so that we can
521 run with multiple seeds.

522 **3. Theory assumptions and proofs**

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

525 Answer: [Yes]

526 Justification: To the best of our knowledge, we clearly state all assumptions. All proof
527 sketches in the main paper are accompanied by deferred full proofs in Appendix A. To the
528 best of our knowledge, we have included references to any theorems or references that the
529 proofs rely on.

530 **4. Experimental result reproducibility**

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

534 Answer: [Yes]

535 Justification: We provide the explicit high-level algorithm, as well as recommendations for
536 low-level details. We also provide hyperparameters toward the goal of reproducibility.

537 **5. Open access to data and code**

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

541 Answer: [Yes]

542 Justification: We provide the code here: [https://anonymous.4open.science/
543 r/0E1FjO/README.md](https://anonymous.4open.science/r/0E1FjO/README.md).

544 **6. Experimental setting/details**

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

547 Answer: [Yes]

548 Justification: We provide experimental settings in the paper when the experiments are
549 introduced. We do a best-faith effort to include the details relevant for appreciating the
550 results. We provide full experimental settings in Section B.

551 **7. Experiment statistical significance**

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

554 Answer: [Yes]
555 Justification: We report error bars when there are multiple seeds.

556 **8. Experiments compute resources**

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

560 Answer: [Yes]
561 Justification: We provide compute details.

562 **9. Code of ethics**

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

565 Answer: [Yes]
566 Justification: It does.

567 **10. Broader impacts**

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

570 Answer: [N/A]
571 Justification: The paper is foundational research and not tied to particular applications.

572 **11. Safeguards**

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

576 Answer: [N/A]
577 Justification: The paper poses no such risks.

578 **12. Licenses for existing assets**

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

582 Answer: [Yes]
583 Justification: We do a best-faith effort to cite all previous assets, and we are happy to include
584 any references that we may have overlooked.

585 **13. New assets**

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

588 Answer: [Yes]
589 Justification: We include as much detail as possible about our code.

590 **14. Crowdsourcing and research with human subjects**

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

594 Answer: [N/A]
595 Justification: The paper does not involve crowdsourcing nor research with human subjects.

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

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

602 Answer: [N/A]

603 Justification: The paper does not involve crowdsourcing nor research with human subjects.

604 **16. Declaration of LLM usage**

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

609 Answer: [N/A]

610 Justification: The core method development in this research does not involve LLMs as any
611 important, original, or non-standard components.