

CoEvo-RL: Co-Evolving Policies and Experiential Memories for LLM Agents

Anonymous ACL submission

Abstract

Experience-augmented LLM agents can learn from both parametric reinforcement learning and explicit trajectory-level memories, but existing methods often optimize the policy and the experience bank in isolation. This separation creates a mismatch: as the policy improves, earlier experiences may be internalized and lose marginal value, while a stale or independently updated bank may fail to provide the knowledge the current policy actually needs. We propose CoEvo-RL, a policy-experience co-evolution framework that updates a GRPO-trained policy and an experiential memory bank concurrently. Across AppWorld, τ^2 -Bench, and LifelongAgentBench, CoEvo-RL outperforms policy-only RL, experience-only learning, and other experience-augmented baselines, improving average success rate from 0.781 to 0.820. Further analyses show faster learning and persistent gains from refreshed experience, indicating that co-evolution improves both final performance and training dynamics.

1 Introduction

LLM agents need learning mechanisms beyond a single policy update (Zhang et al., 2026b). Interactive LLM agents must repeatedly perceive an environment, select actions, observe feedback, and improve over long horizons. Reinforcement learning can improve the parametric policy from task rewards, but many agent failures are also episodic and reusable: a failed trajectory may reveal a missing precondition, a tool-use convention, or a domain-specific action pattern that should be preserved explicitly. Recent experience-augmented agents therefore store reflections, skills, or distilled lessons from past trajectories in an external experiential memory bank and retrieve them as context for future decisions (Shinn et al., 2023; Wang et al., 2023; Zhao et al., 2024).

Optimizing policy and experience in isolation is limiting. Existing learning pipelines typically

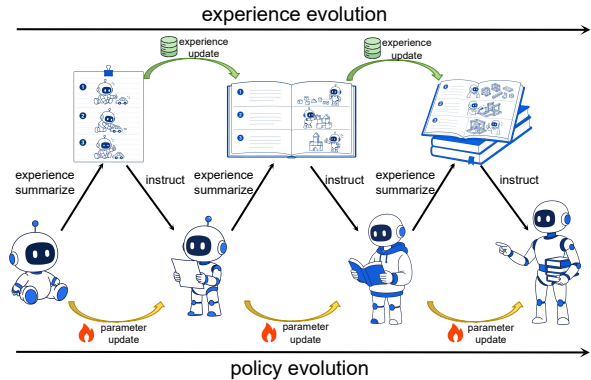


Figure 1: The key idea of CoEvo-RL. A stronger parametric policy can discover more useful experiences, while better experiences guide the policy to complete tasks more effectively.

emphasize one side of this system. A policy-only agent can improve its model parameters from environment rewards, but it lacks a persistent non-parametric store for preserving explicit lessons about failures, partial successes, and task-specific decision conditions. An experience-only agent can accumulate reusable entries, but a frozen policy cannot internalize recurring action patterns or learn when and how to exploit retrieved experience. This mismatch becomes especially important in sparse-reward agent benchmarks, where early rollouts are often uniformly unsuccessful and provide weak advantage signals for policy learning. Conversely, useful experience entries are hard to write before the policy has discovered informative trajectories.

In this paper, we propose CoEvo-RL, which moves beyond isolated policy or experience optimization by jointly learning a parametric policy and a non-parametric experiential memory bank. As shown in Figure 1, the core idea is that policy improvement and experience construction should reinforce each other: a stronger policy produces more informative trajectories for experience extraction, while better experience entries reshape future

067 rollouts and provide more useful learning signals
068 for policy optimization.

069 COEVO-RL alternates between two update steps
070 over a shared stream of agent trajectories. In the
071 policy step, the agent retrieves relevant entries from
072 the current experiential memory bank and performs
073 experience-conditioned rollouts; the resulting re-
074 wards are then used to update the parametric policy
075 with GRPO. In the experience step, the same tra-
076 jectories are converted into candidate experience
077 entries that summarize reusable strategies, task con-
078 straints, or failure patterns. Candidate entries are
079 committed only when replay indicates that they im-
080 prove over the current experiential memory bank.
081 This creates a closed training loop: the experiential
082 memory bank shapes the rollout distribution seen
083 by the policy, while the improved policy produces
084 new trajectories from which better entries can be
085 extracted.

086 We empirically evaluate COEVO-RL on three in-
087 teractive agent benchmarks: AppWorld, τ^2 -Bench,
088 and LifelongAgentBench. Compared with policy-
089 only RL, experience-only learning, and prior
090 policy-experience training, COEVO-RL achieves
091 the best success rate across all evaluated settings.
092 Beyond final performance, our ablations and train-
093 ing dynamics show why co-evolution matters: an
094 evolving bank continues to expose policy-relevant
095 knowledge after earlier entries have been partly
096 internalized.

097 Our main contributions are:

- 098 • We propose CoEvo-RL, a policy–experience co-
099 evolving pipeline that alternates between policy
100 updates and experience extraction, converting
101 improved agent behaviors into reusable entries
102 while using accumulated experience to further
103 guide subsequent policy learning.
- 104 • We evaluate CoEvo-RL on AppWorld, τ^2 -Bench,
105 and LifelongAgentBench, comparing against
106 GRPO, experience-only pipelines (MemRL), and
107 existing policy-experience baselines (SkillRL).
108 The results demonstrate that our method achieves
109 more effective optimization and consistently out-
110 performs competing approaches.
- 111 • We provide ablations and training-dynamics anal-
112 yses showing that policy learning and experience
113 evolution are mutually reinforcing, and that con-
114 tinually refreshed experience are necessary as
115 earlier entries are gradually internalized by the
116 policy.

2 Related Work 117

2.1 Reinforcement Learning for LLM Agents 118

119 Reinforcement learning has become a central ap-
120 proach for improving LLM agents on interactive
121 tasks, where agents must learn effective action poli-
122 cies from sparse task feedback. Recent methods
123 often adopt critic-free group-relative policy opti-
124 mization (GRPO) to improve training stability un-
125 der verifiable or outcome-based rewards. Further-
126 more, GiGPO extends this idea to agent tasks by
127 incorporating anchor-state step-level advantages
128 in addition to episode-level feedback (Feng et al.,
129 2025); and RLVMR combines outcome rewards
130 with verifiable meta-reasoning rewards to improve
131 generalization on embodied text-based environ-
132 ments (Zhang et al., 2025). These methods demon-
133 strate the effectiveness of RL for optimizing the
134 parametric policy of LLM agents. However, they
135 typically treat experiential memory as fixed compo-
136 nents rather than optimization targets. In contrast,
137 our work targets a different direction: we treat both
138 the parametric policy and experiential memory as
139 learnable components, and optimize them jointly
140 under task reward.

2.2 Memory Systems for LLM Agents 141

142 Experiential memory has emerged as an important
143 direction for LLM agents, aiming to distill reusable
144 knowledge from the agent’s own past trajectories
145 and reuse it on future tasks. Most existing systems
146 build such memories on top of a frozen policy. For
147 example, Reflexion stores verbal self-critique from
148 failed trajectories (Shinn et al., 2023), Voyager ac-
149 cumulates executable skills (Wang et al., 2023),
150 ExpeL distills cross-task natural-language insights
151 (Zhao et al., 2024), A-Mem organizes intercon-
152 nected memory notes (Xu et al., 2025), and Mem0
153 consolidates long interaction histories into hybrid
154 vector-graph memory (Chhikara et al., 2025). Be-
155 yond these frozen-policy memory pipelines, a more
156 recent line makes memory construction itself learn-
157 able, such as MemSkill (Zhang et al., 2026a),
158 EvoSkill (Alzubi et al., 2026), and AtomMem (Huo
159 et al., 2026), which learn how to update or orga-
160 nize reusable memories. However, these methods
161 still mainly optimize the memory mechanism it-
162 self, while the task policy is either frozen or not
163 jointly evolved with the memory bank. In contrast,
164 our CoEvo-RL framework treats policy learning
165 and memory optimization as a coupled process:
166 reward-driven policy updates produce improved

trajectories for memory extraction, and the resulting memories further support subsequent policy learning.

3 Method

We present COEVO-RL, a framework that jointly evolves a parametric policy and a non-parametric experience bank from the same stream of agent trajectories. The policy learns from experience-conditioned rollouts, while new experience entries are distilled from trajectories produced by stronger policies. Figure 2 gives an overview of this shared training loop and its two coupled update paths.

3.1 Preliminaries

We consider an LLM agent interacting with an environment on tasks sampled from a distribution \mathcal{D} . A task is specified by a natural-language instruction $x \sim \mathcal{D}$. During execution, the agent produces a multi-turn trajectory $\tau = (a_1, o_1, \dots, a_T, o_T)$ and receives a terminal reward $R(\tau) \in [0, 1]$ from the environment.

At training step t , the agent state consists of two learnable components:

$$S_t = (\theta_t, E_t), \quad (1)$$

where θ_t parameterizes the policy π_{θ_t} and $E_t = \{e_1, e_2, \dots\}$ is an experience bank of distilled trajectory-level knowledge. Each experience entry stores a reusable lesson from previous trajectories, including what strategy to apply and under what task condition it is useful. Given a task x , a query $q(x)$ is derived from the task instruction and scored against each experience entry with a scoring function s :

$$C_t(x) = \text{TopK}_{e \in E_t}(s(q(x), e), k). \quad (2)$$

where k is the retrieval budget and $C_t(x)$ is the compact experience block inserted into the rollout prompt. In our implementation, s is instantiated with structured TF-IDF similarity. The policy then acts under the experience-conditioned context:

$$\tau \sim \pi_{\theta_t}(\cdot \mid x, C_t(x)). \quad (3)$$

COEVO-RL updates both components of S_t through a co-evolutionary loop: policy optimization changes how trajectories are generated, and experience evolution changes what reusable knowledge conditions future rollouts. The two updates are coupled through a shared rollout stream.

3.2 Policy-Experience Co-Evolution

COEVO-RL alternates between policy and experience updates on a shared rollout stream. At each policy update, the current bank E_t conditions rollouts through retrieved experiences, providing task-specific constraints, failure warnings, and action recipes; the resulting rewards update θ_t with policy optimization.

Periodically, recent trajectories from the updated policy are distilled into candidate experience entries. Candidates are admitted only when they improve future rollouts under the validation procedure in Section 3.4. Thus, retrieved experience shapes the policy’s training distribution, while the improving policy supplies new evidence for experience evolution.

The resulting feedback loop can be summarized as

$$E_t \rightarrow \{\tau_t, R(\tau_t)\} \rightarrow \theta_{t+1} \rightarrow \{\tau_{t+1}\} \rightarrow E_{t+1}. \quad (4)$$

This loop is not a one-time experience initialization step. It runs throughout training, so the policy and experience bank continually adapt to each other’s current capabilities. Algorithm 1 summarizes the resulting training procedure.

3.3 Experience-Conditioned Policy Optimization

For each task x , we retrieve $C_t(x)$ from the current experience bank and sample a group of N trajectories under the same retrieved context:

$$\{\tau^{(i)}\}_{i=1}^N \sim \pi_{\theta_t}(\cdot \mid x, C_t(x)). \quad (5)$$

We optimize the policy with GRPO (Shao et al., 2024). For each trajectory in the group, the advantage is computed from relative task outcomes:

$$A^{(i)} = \frac{R(\tau^{(i)}) - \text{mean}_j R(\tau^{(j)})}{\text{std}_j R(\tau^{(j)}) + \varepsilon}. \quad (6)$$

The normalized advantages are used in the standard clipped GRPO objective with KL regularization to a reference policy.

Since all trajectories in a group share the same $C_t(x)$, their relative advantages compare policy behavior under a fixed experience context. This reduces retrieval-induced variance and leaves experience quality to the experience-evolution step.

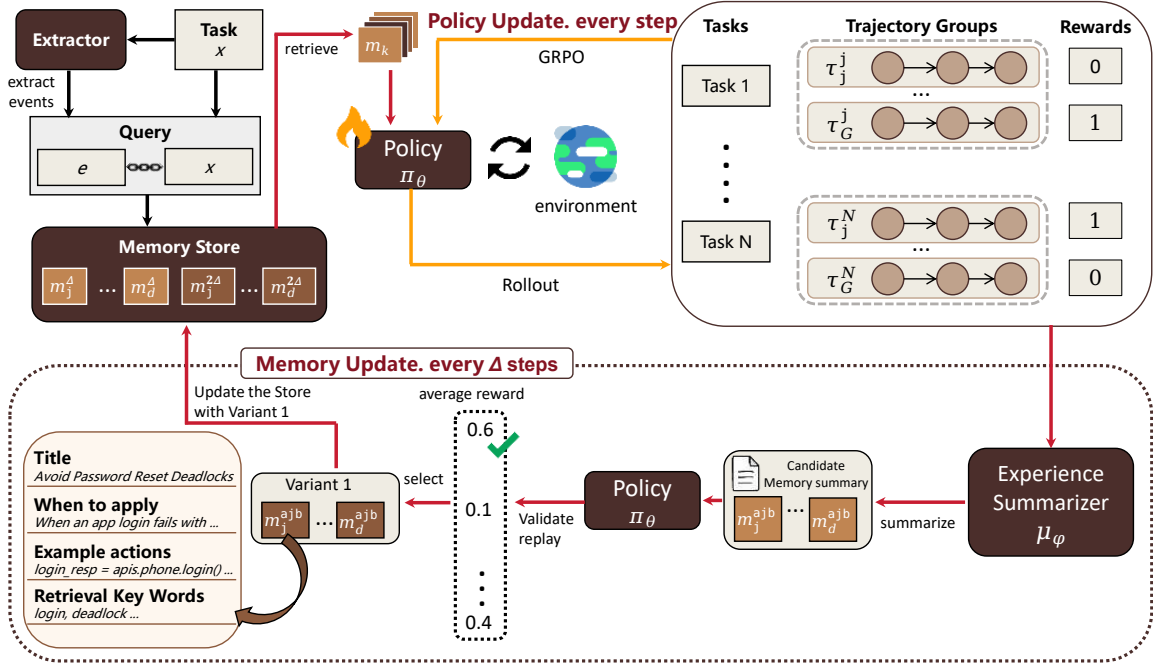


Figure 2: Overview of COEVO-RL, where experience-conditioned rollouts update the policy and the resulting trajectories are distilled into validated experience entries.

3.4 Experience Evolution

Every Δ policy-update steps, we update the experience bank from recent trajectories. The experience updater μ_ϕ receives a small trajectory buffer together with the current bank E_t and distills candidate entries that capture reusable experience. It prioritizes tasks with both successful and failed rollouts, because these contrastive pairs expose the pivotal decision point under the same task context. When such pairs are unavailable, it falls back to success-only clusters for reusable action recipes and failure-only clusters for error causes and corrections.

Each candidate experience entry is written in a structured natural-language form. Besides a reusable strategy and its applicability condition, an entry also includes an example action that grounds the strategy in concrete behavior, as well as retrieval keywords used by the structured retriever. This schema keeps entries readable to the policy while giving the retriever explicit cues for matching future tasks to the right experience.

Because trajectory-derived entries can be noisy, we validate them by replay before committing them to the live bank. For a candidate set \mathcal{V} of new entries, we form $E' = E_t \cup \mathcal{V}$ and compare its mean reward $\hat{R}(E')$ with the current bank $\hat{R}(E_t)$ on the tasks that produced the candidates. The next

experience bank is

$$E_{t+1} = \begin{cases} E', & \text{if } \hat{R}(E') > \hat{R}(E_t) + \delta, \\ E_t, & \text{otherwise.} \end{cases} \quad (7)$$

Thus, the evolving experience bank is tied to task performance instead of treating experience writing as a separate objective.

This design lets the two components co-evolve: few early successes or revealing failures can provide experiences that guide later exploration, and as the policy improves, its trajectories provide stronger evidence for updating the experience bank.

4 Experiments

Our experiments evaluate whether policy learning and experience construction are more effective when they evolve together rather than in isolation. We organize the evaluation around four questions: (1) Does co-evolution improve over policy-only, experience-only, and existing policy-experience baselines? (2) Which retrieval mechanism best selects useful experience for the current policy? (3) Does the timing of experience evolution matter? (4) Do the learning dynamics show evidence that experience is progressively internalized while the bank keeps exposing new policy-relevant knowledge?

Algorithm 1: Policy-Experience Co-Evolution in COEVO-RL

Input: Initial policy π_{θ_0} , bank E_0 , updater μ_ϕ , interval Δ , budget k

Output: Optimized policy π_θ and evolved bank E

```
1 for  $t = 0, 1, 2, \dots$  do
2   Sample a task batch  $\mathcal{B}_t \sim \mathcal{D}$ ;
3   foreach  $x \in \mathcal{B}_t$  do
4      $C_t(x) \leftarrow \text{TopK}_{e \in E_t}(s(q(x), e), k)$ ;
5      $\{\tau_x^{(i)}\}_{i=1}^N \sim \pi_{\theta_t}(\cdot | x, C_t(x))$ ;
6      $\theta_{t+1} \leftarrow \text{GRPO}(\theta_t, \{(\tau, R(\tau))\})$ ;
7     if  $t \bmod \Delta = 0$  then
8        $\mathcal{V} \leftarrow \mu_\phi(\{\tau\}, E_t)$ ;
9        $E^+ \leftarrow E_t \cup \mathcal{V}$ ;
10      Estimate  $\hat{R}(E^+)$  and  $\hat{R}(E_t)$  by
          replay;
11      if  $\hat{R}(E^+) > \hat{R}(E_t) + \delta$  then
12         $E_{t+1} \leftarrow E^+$ ;
13      else
14         $E_{t+1} \leftarrow E_t$ ;
15    else
16       $E_{t+1} \leftarrow E_t$ ;
```

4.1 Experimental Setup

Benchmarks. We evaluate on three interactive agent benchmarks. **AppWorld** (Trivedi et al., 2024) evaluates interactive coding agents in a simulated world of everyday applications. τ^2 -**Bench** (Yao et al., 2024) evaluates agents in customer-service workflows that require following domain policies; we use its *telecom* domain following the Artificial Analysis intelligence-benchmarking methodology¹. **LifelongAgent-Bench** (Zheng et al., 2025) (LLB) provides skill-grounded interactive tasks in three environments: database, operating system, and knowledge graph. Following MemRL (Zhang et al., 2026b), we use the DBBENCH database domain and the OSINTERACTION operating-system domain with the same train/test split as MemRL.

Baselines and Metrics. We compare our method against three method families: **Policy-only** (GRPO with no experience module), **Experience-only** (MemRL (Zhang et al., 2026b)), and **Policy &**

¹<https://artificialanalysis.ai/methodology/intelligence-benchmarking>

Experience (SkillRL (Xia et al., 2026)), which update both policy and experience during training. All benchmarks are evaluated by task-level *success rate*, using the official success criterion of each benchmark.

Implementation Details. In our experiments, all methods share the same policy backbone (QWEN3-8B (Yang et al., 2025)) and are trained with the GRPO update of Eq. (6) ($N=8$ rollouts per task, batch of 16 tasks, learning rate 10^{-6} , low-variance KL coefficient 10^{-3}) on a single $4 \times A100$ -80G node with asynchronous SGLang rollout. The experience pipeline of §3.4 runs every Δ steps ($\Delta=20$ on AppWorld/LLB, $\Delta=15$ on τ^2 -Bench) with retrieval top- $k=8$; the default retriever uses the structured TF-IDF scoring of Eq. (2) and the experience updater μ_ϕ is a closed-source LLM API.

4.2 Main Results

Table 1 summarizes the main comparison across four task settings that cover tool use, agentic workflow following, operating-system interaction, and database operation. CoEvo-RL achieves the best result on every benchmark and improves the average success rate from 0.781 with SkillRL to 0.820. We draw three observations.

First, policy optimization and experience evolution are complementary. Compared with the policy-only GRPO baseline, CoEvo-RL improves the average success rate by 8.0 points (0.740 \rightarrow 0.820). In contrast, experience-only methods are weaker than policy optimization on average: MemRL reaches 0.606, and *Ours-Exp* reaches 0.571 under a frozen policy. External experience therefore complements, rather than replaces, policy optimization.

Second, reward-validated experience evolution improves over prior policy-experience training. SkillRL also augments reinforcement learning with evolving skills distilled from the current policy’s rollouts, making it a strong comparison for coupled policy-experience learning. Yet CoEvo-RL improves over SkillRL on all four settings: +5.4 points on AppWorld, +4.5 points on τ^2 -Telecom, +2.0 points on LLB-OS, and +3.6 points on LLB-DB. This consistent gain suggests that effective policy-experience co-evolution requires not only shared on-policy rollouts, but also a shared optimization objective, with environment reward guiding which experiences are retained.

Third, the advantage is largest on environments where explicit reusable constraints matter. Com-

Family	Method	AppWorld	τ^2 -Telecom	LLB-OS	LLB-DB	Avg.
None	Qwen3-8B (frozen)	0.211	0.225	0.493	0.791	0.430
Policy-only	GRPO (Shao et al., 2024)	0.478	0.900	0.667	0.914	0.740
Experience-only	MemRL (Zhang et al., 2026b)	0.481	0.475	0.626	0.840	0.606
	Ours-Exp (frozen policy)	0.454	0.350	0.653	0.827	0.571
Policy & Experience	SkillRL (Xia et al., 2026)	0.554	0.930	0.700	0.942	0.781
	Ours	0.608	0.975	0.720	0.978	0.820

Table 1: Main success-rate results across AppWorld, τ^2 -Bench, and two LifelongAgentBench (LLB) settings. The average is computed over the four task settings.

pared with LLB-OS and LLB-DB, whose domains are already familiar to the base model, AppWorld and τ^2 -Bench telecom impose more environment-specific interaction rules, hidden APIs, and procedural manuals. On difficult tasks in these environments, policy rollouts alone rarely produce successful trajectories under sparse terminal rewards, so RL receives too little positive signal to reliably learn the required rules from scratch. These are exactly the settings where CoEvo-RL gains most over SkillRL and GRPO, suggesting that co-evolved experience is most useful when the base policy is not familiar with the environment.

4.3 Ablation Studies

We next isolate the design choices in CoEvo-RL: how entries are retrieved, how often the experience bank is updated, how policy and experience updates are scheduled, and how many entries should be placed in the context. Unless otherwise specified, ablations use the same backbone and training recipe as §4.1.

4.3.1 Effectiveness of Retrieval Methods

We evaluate different retrievers in Eq. (2) using 47 query–target pairs, where each experience’s recorded source task serves as the query and the experience itself is the ground-truth target. We report Hit@1, Hit@3, MRR, and downstream success rate.

Table 2 shows that TF-IDF achieves the best retrieval quality and downstream performance, with 88.6% Hit@1, 95.5% Hit@3, 0.871 MRR, and a 72.0% success rate. In contrast, embedding similarity has the lowest Hit@1 (65.9%), MRR (0.754), and success rate (65.3%), indicating that broad semantic similarity alone misses task-specific operational cues such as command names, API objects, file paths, and retrieval keys. We therefore use

Retriever	Hit@1	Hit@3	MRR	SR (%)
Embedding similarity	65.9	81.8	0.754	65.3
TF-IDF (Ours)	88.6	95.5	0.871	72.0
BM25	81.8	93.2	0.849	68.0
TF-IDF + Embedding	79.5	89.7	0.861	68.7
BM25 + Embedding	75.0	87.4	0.865	67.3

Table 2: Comparison of retrievers on LLB-OSInteraction. A retrieval is judged correct if the retrieved entry is distilled from the same task as the current task. **Success Rate** is the downstream success rate on the test split under each retriever.

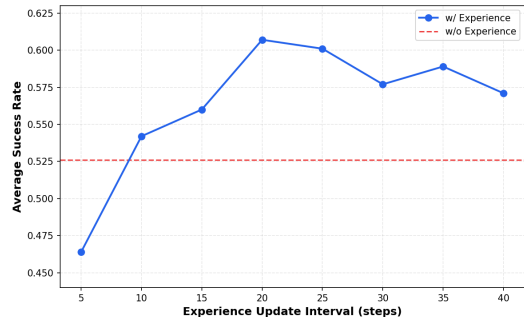


Figure 3: Effect of the experience update interval Δ on AppWorld. One AppWorld epoch corresponds to 5 training steps; we sweep $\Delta \in \{5, 10, 15, 20, 25, 30, 35, 40\}$ steps (i.e., 1–8 epochs). The y -axis is task success rate at convergence; all other hyper-parameters follow §4.1.

TF-IDF as the default retriever in subsequent experiments.

4.3.2 Sensitivity to Update Interval

The experience update interval controls a trade-off between freshness and stability. Updating too frequently can commit entries from noisy early rollouts, while updating too sparsely leaves the policy with stale experience. We therefore sweep Δ from 5 to 40 steps on AppWorld, where 5 steps correspond to one training epoch.

Figure 3 shows a clear non-monotonic relation-

Training schedule	Success Rate
P→E	0.68
E→P	0.70
Joint Co-Evolution (Ours)	0.72

Table 3: Effect of the policy-experience training schedule on LIFELONGAGENTBENCH OSInteraction.

ship between update interval and final performance. The success rate rises from 0.464 at $\Delta=5$ to a peak of 0.607 at $\Delta=20$, then gradually drops as updates become less frequent. Updating every epoch is even worse than the GRPO-without-experience reference (0.526). A likely explanation is twofold. First, very early rollouts yield immature entries that are not yet reliably reusable. Second, if the bank is refreshed too often, the retrieved context distribution seen by the policy changes too quickly across training, increasing non-stability in the learning signal and making credit assignment noisier. At the other extreme, overly infrequent updates leave the bank stale. The best interval therefore balances freshness and stability, allowing useful experience to remain available long enough for the policy to benefit from it.

4.3.3 Effect of Co-Evolution Schedule

We compare three policy-experience training schedules on LIFELONGAGENTBENCH OSInteraction. In P→E, the policy is first trained without experience, then frozen while the experience bank evolves. In E→P, the experience bank is first grown under a frozen policy, then fixed while the policy is trained. These staged schedules isolate the two learning channels, whereas our joint co-evolution schedule alternates policy and experience updates every Δ steps under the same total compute budget. This comparison tests whether policy and experience benefit from adapting to each other during training.

As shown in Table 3, joint co-evolution achieves the highest success rate (0.72), outperforming both staged alternatives. The results suggest that staged schedules may introduce a mismatch between one component and the distribution induced by the other, whereas alternating updates better maintain alignment between retrieved entries and the current policy’s failure modes.

4.3.4 Effect of Retrieved Experience Num.

The retrieval budget controls the amount of external experience available at inference time, trading

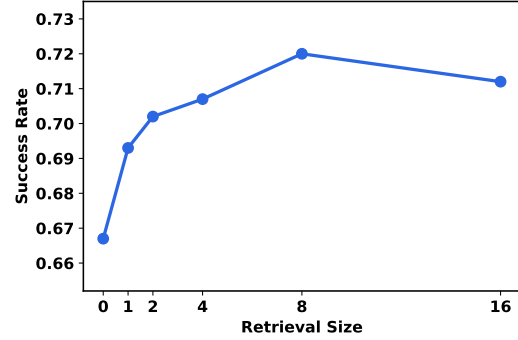


Figure 4: Effect of the budget of retrieved experience entries on LLB-OSInteraction. x -axis: number of retrieved entries; y -axis: success rate.

broader coverage against longer contexts and potentially distracting redundancy. To examine this effect, we evaluate CoEvo-RL on LLB-OSInteraction with varying entry budgets while keeping the retriever and policy fixed.

Figure 4 shows that performance improves as the budget increases from 0 to 8 entries (0.667 \rightarrow 0.720), then slightly drops at 16 entries (0.712). This suggests that retrieval benefits come from a small set of high-precision entries rather than from exposing the policy to as much external text as possible. We therefore use a relatively compact retrieval budget in the default setting.

5 Discussion

We further analyze how co-evolving experience shapes CoEvo-RL’s training dynamics.

5.1 Convergence and Sample Efficiency

Final success rates summarize the endpoint of training, but they do not characterize how quickly performance improves. We therefore use learning curves to examine whether co-evolving experience changes the learning dynamics. If validated entries provide reusable task constraints during training, CoEvo-RL should begin to improve earlier than policy-only optimization rather than matching it until the final stage.

Figure 5 shows this pattern across all three benchmarks. CoEvo-RL rises earlier than the policy-only baseline and generally maintains a higher success rate throughout training, with especially clear gains before convergence on AppWorld and τ^2 -Bench telecom. This suggests that experience updates improve sample efficiency by exposing the policy to validated constraints and strategies.

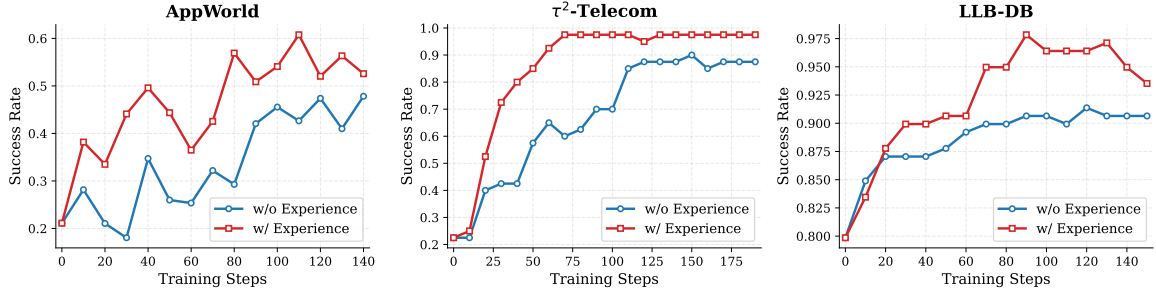


Figure 5: Learning curves on AppWorld, τ^2 -Bench telecom, and LLB-DB. x -axis: training step; y -axis: success rate. We compare *Policy-only* (GRPO without retrieved experience) and *CoEvo-RL* (joint policy-experience co-evolution).

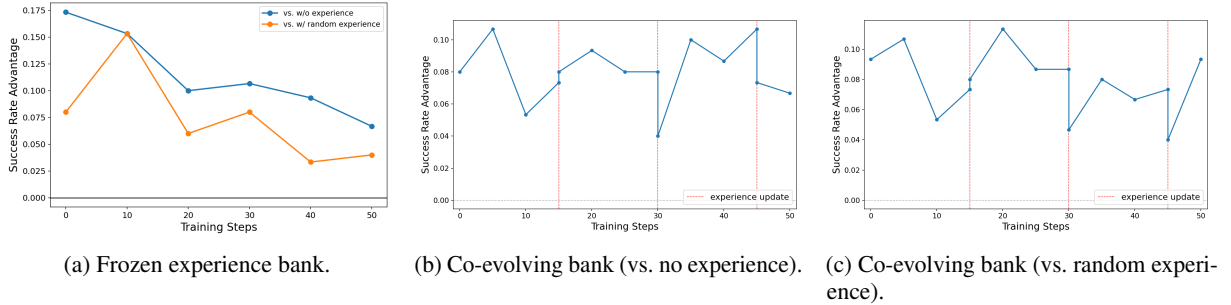


Figure 6: Experience advantage dynamics on LIFELONGAGENTBENCH OSInteraction. x -axis: training step; y -axis: success-rate advantage of top-1 retrieved experience over each baseline. (a) Initial experience bank is frozen. (b–c) Experience is updated every 15 steps.

5.2 Experience Internalization Dynamics

We study whether the policy progressively internalizes retrieved experience on LIFELONGAGENTBENCH OSInteraction. At each checkpoint θ_t , we compare top-1 retrieved experience with two baselines that either remove experience or inject one random entry from the same bank. We use top-1 rather than the top- k setting in the main experiments to sharpen credit assignment, since a larger k makes the random baseline more likely to include a useful entry. Figure 6 reports the success-rate advantage under a frozen initial bank and a co-evolving bank updated every 15 steps.

With a frozen bank, the advantage steadily shrinks, falling from about 0.175 to 0.075 against no experience and from an early peak of about 0.15 to roughly 0.04 against random experience. This decline suggests that the policy gradually absorbs the behavior encoded in the initial entries. With a co-evolving bank, the advantage remains recurrent instead of vanishing, staying around 0.04–0.11 against no experience and 0.04–0.12 against random experience. It also dips just after updates, most visibly near step 30, where the advantage falls to about 0.04–0.05 before recovering. These dips suggest that newly committed entries capture er-

rors and decisions that the current policy has not yet learned to exploit, so their value appears after several policy updates. Thus, policy-experience co-evolution is necessary because internalized entries lose marginal value, while the bank must keep introducing new experience the current policy has not yet absorbed.

6 Conclusion

We study reinforcement learning for experience-augmented LLM agents from a policy-experience co-evolution perspective. COEVO-RL couples GRPO-based policy optimization with non-parametric experience construction through shared rollouts, where structured TF-IDF retrieval provides relevant experience and periodic replay validation commits only reward-improving entries. Experiments on AppWorld, τ^2 -Bench, and LifelongAgentBench show that jointly optimizing policy and experience consistently outperforms optimizing either component alone. Further analyses validate the effects of retrieval design, update schedules, proposal validation, convergence behavior, and experience internalization. Future work may explore cross-environment experience transfer and compact experience-bank maintenance.

553
554
555
556
557
558
559
560

561

562
563
564
565

566
567
568
569

570
571
572

573
574
575
576

577
578
579
580
581
582

583
584
585
586
587

588
589
590
591
592
593
594
595

596
597
598
599
600

601
602
603
604

Limitations

The experience updater is instantiated with a closed-source LLM API, so its parameters cannot be updated to summarize entries that best match the evolving needs of the policy. We also leave experience consolidation for future work, including how to remove stale internalized entries and prevent unbounded growth of the experience bank.

References

Salaheddin Alzubi, Noah Provenzano, Jaydon Bingham, Weiyuan Chen, and Tu Vu. 2026. EvoSkill: Automated skill discovery for multi-agent systems. *arXiv preprint arXiv:2603.02766*.

Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet Singh, and Deshraj Yadav. 2025. Mem0: Building production-ready AI agents with scalable long-term memory. *arXiv preprint arXiv:2504.19413*.

Lang Feng, Zhenghai Xue, Tingcong Liu, and Bo An. 2025. Group-in-group policy optimization for LLM agent training. *arXiv preprint arXiv:2505.10978*.

Yupeng Huo, Yaxi Lu, Zhong Zhang, Haotian Chen, and Yankai Lin. 2026. *Atomem: Learnable dynamic agentic memory with atomic memory operation*. Preprint, arXiv:2601.08323.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.

Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. In *Advances in Neural Information Processing Systems*.

Harsh Trivedi, Tushar Khot, Mareike Hartmann, Ruskin Manku, Vinty Dong, Edward Li, Shashank Gupta, Ashish Sabharwal, and Niranjan Balasubramanian. 2024. Appworld: A controllable world of apps and people for benchmarking interactive coding agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16022–16076.

Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023. Voyager: An open-ended embodied agent with large language models. In *Advances in Neural Information Processing Systems*.

Peng Xia, Jianwen Chen, Hanyang Wang, Jiaqi Liu, Kaide Zeng, Yu Wang, Siwei Han, Yiyang Zhou, Xujiang Zhao, Haifeng Chen, Zeyu Zheng, Cihang Xie, and Huaxiu Yao. 2026. Skillrl: Evolving agents

via recursive skill-augmented reinforcement learning. Preprint, arXiv:2602.08234. 605
606

Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. 2025. A-Mem: Agentic memory for LLM agents. In *Advances in Neural Information Processing Systems*. 607
608
609
610

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*. 611
612
613
614
615

Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024. τ -bench: A benchmark for tool-agent-user interaction in real-world domains. *arXiv preprint arXiv:2406.12045*. 616
617
618
619

Haozhen Zhang, Quanyu Long, Jianzhu Bao, Tao Feng, Weizhi Zhang, Haodong Yue, and Wenya Wang. 2026a. MemSkill: Learning and evolving memory skills for self-evolving agents. *arXiv preprint arXiv:2602.02474*. 620
621
622
623
624

Shengtao Zhang, Jiaqian Wang, Ruiwen Zhou, Junwei Liao, Yuchen Feng, Zhuo Li, Yujie Zheng, Weinan Zhang, Ying Wen, Zhiyu Li, and 1 others. 2026b. Memrl: Self-evolving agents via runtime reinforcement learning on episodic memory. *arXiv preprint arXiv:2601.03192*. 625
626
627
628
629
630

Zijing Zhang, Ziyang Chen, Mingxiao Li, Zhaopeng Tu, and Xiaolong Li. 2025. RLVMR: Reinforcement learning with verifiable meta-reasoning rewards for robust long-horizon agents. *arXiv preprint arXiv:2507.22844*. 631
632
633
634
635

Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. 2024. Expel: Llm agents are experiential learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 636
637
638
639

Junhao Zheng, Xidi Cai, Qiuke Li, Duzhen Zhang, Zhongzhi Li, Yingying Zhang, Le Song, and Qianli Ma. 2025. Lifelongagentbench: Evaluating llm agents as lifelong learners. *arXiv preprint arXiv:2505.11942*. 640
641
642
643
644