

---

# FlashEvolve: Accelerating Agent Evolution with Asynchronous Stage Orchestration

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 LLM-based evolution has emerged as a promising way to improve agents by re-  
2 fining non-parametric artifacts, but its wall-clock cost remains a major bottleneck.  
3 We identify that this cost comes from synchronized stage execution and imbalance  
4 inside each LLM-heavy stage. We present FlashEvolve, an efficient framework that  
5 replaces synchronized execution with asynchronous workers and queues, allowing  
6 different stages and steps to overlap. To handle data staleness introduced by asyn-  
7 chrony, FlashEvolve tracks artifact versions and applies different policies to update,  
8 discard, or patch stale artifacts. Unlike weight-space staleness in asynchronous RL,  
9 language-space staleness is inspectable and repairable: a stale artifact is not just  
10 delayed work, but readable evidence that the LLM can reflect on, revise, and turn  
11 into useful evolution signal. FlashEvolve further improves throughput and token  
12 efficiency with speculative stage completion and adaptive workflow control.

13 On GEPA workloads, FlashEvolve improves proposal throughput by  $3.5\times$  on local  
14 vLLM and  $4.9\times$  on API serving over synchronous GEPA. The same design also  
15 applies to ACE and Meta-Harness. Our code repository is available at <https://anonymous.4open.science/r/FlashEvolve-FEC7/>.  
16

## 17 1 Introduction

18 A growing line of recent work enables LLM agents to evolve themselves. Instead of updating model  
19 weights, these systems iteratively refine the non-parametric components that govern their behavior,  
20 including system prompts [2, 28, 29], context and memory [34, 22, 33], harness code [16, 14] and  
21 generated programs [20, 13, 3]. This emerging paradigm of test-time self-evolution [6] fundamentally  
22 relaxes the access requirements of weight-space adaptation: it requires neither the labeled trajectories  
23 used by supervised fine-tuning nor the gradient updates required by reinforcement learning. By having  
24 an LLM reflect on full execution traces rather than optimize against scalar rewards, this paradigm  
25 draws a richer learning signal from each rollout: GEPA [2] outperforms GRPO with an average gain  
26 of 6% across six reasoning benchmarks, while Meta-Harness [14] automatically discovers agent  
27 harnesses that surpass the best hand-engineered baselines on different domain-specific benchmarks.

28 Despite its algorithmic appeal, agent evolution remains expensive in wall-clock execution time.  
29 Existing evolution algorithms pursue “faster” evolution by improving the quality of each step through  
30 stronger reflection [2, 34, 32], better artifact proposal and search [13, 18], or larger-batch updates [15],  
31 thereby reducing the number of steps needed. However, fewer evolution steps do not necessarily  
32 translate into shorter wall-clock time. For example, on IFBench, a single GEPA evolution step already  
33 takes  $\sim 2$  minutes; Combee [15] parallelizes proposal generation, but further stretches each step to  
34  $\sim 2.8$  minutes. Reaching a stable improvement requires more than 2 hours on an H100 GPU. This  
35 cost further grows with data scale, making evolution runs slow to tune and deploy in practice.

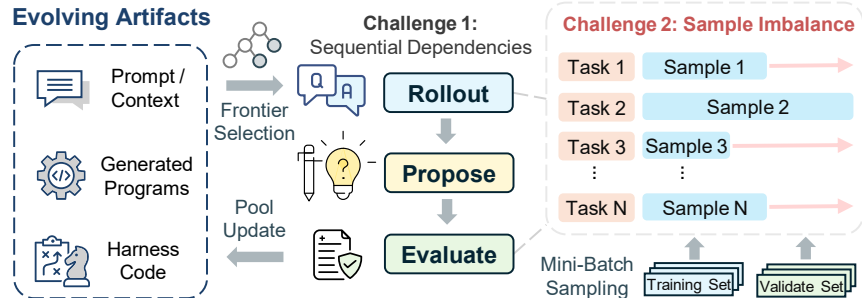


Figure 1: Illustration of the multi-stage execution in agent evolution. The synchronized stage orchestration in existing implementations [2, 34, 15] exposes two efficiency challenges: sequential dependencies across stages and sample workload imbalance within each individual stage.

36 Such high wall-clock cost comes from **synchronized stage execution**. As shown in Figure 1, each  
 37 evolution step runs a sequence of LLM-heavy stages, such as running the current artifact on a mini-  
 38 batch of inputs, proposing a new candidate artifact, and evaluating the new one. A later stage cannot  
 39 start until the previous stage has fully completed. Such serial structure prevents overlap across stages.

40 The cost inefficiency is amplified by **generation imbalance** inside each stage. Request lengths vary  
 41 widely across samples, such as different validation samples in the evaluate stage. This creates a  
 42 long-tail effect: the longest requests determine the execution time of the whole stage. This reduces  
 43 the effective batch size in both local serving frameworks [12, 37] and API-based remote calls, leading  
 44 to low resource utilization and inefficient waiting for long samples.

45 To this end, we present FlashEvolve, a framework that improves the time efficiency of agent evolution  
 46 through asynchronous stage orchestration. FlashEvolve treats an evolution loop as a set of LLM-heavy  
 47 stages connected by queues. This allows artifact execution, proposal generation, evaluation, and pool  
 48 update to overlap in time, turning a synchronized loop into a streaming execution pipeline.

49 This design introduces new systems challenges. Asynchronous execution can generate stale items  
 50 because an artifact pool may change while earlier items are still waiting in queues. FlashEvolve  
 51 handles this with artifact-version tracking and staleness-aware policies, including version comparison  
 52 and discarding, or reflective patching for stale language artifacts. This property is specific to agent  
 53 evolution. Unlike weight updates in SFT or reinforcement learning, evolution artifacts are prompts,  
 54 memories, harness code, or programs. A stale artifact is therefore still an inspectable object: its  
 55 relation to the current pool can be judged as complementary, redundant, or conflicting, and can be  
 56 revised by the same LLM mechanism used for proposal. This makes staleness a semantic repair  
 57 problem rather than only a scheduling hazard. FlashEvolve further reduces waiting inside long stages  
 58 through speculative completion, and uses adaptive workflow control to balance workload across  
 59 stages. Together, these mechanisms improve throughput while preserving the quality of evolution.

## 60 2 Background and Motivation

### 61 2.1 Agent Evolution: Self-Improvement Beyond Weight Updates

62 Agent evolve has emerged as a new paradigm for adapting LLM-based systems to new data and  
 63 tasks [6, 4]. This success stems from the already strong reasoning capability of modern LLMs [7, 10],  
 64 which enables a single model to reflect on its own trajectories [26], critique its own outputs [19], and  
 65 propose new artifacts that govern its own behavior, ranging from prompts, memory, and harness code  
 66 that govern how the agent operates, to generated programs that constitute the task solution. Crucially,  
 67 this happens without modifying model weights, sidestepping the training infrastructure of supervised  
 68 fine-tuning [27, 36] and reinforcement learning [7, 21] while delivering comparable or stronger gains.

69 For example, GEPA [2] and ACE [34] use reflection on execution traces to evolve system prompts  
 70 and contextual playbooks. Meta-Harness [14] and AutoHarness [16] use a coding agent to evolve the  
 71 harness based on prior runs and their failure modes. AlphaEvolve [20] and ShinkaEvolve [13] push

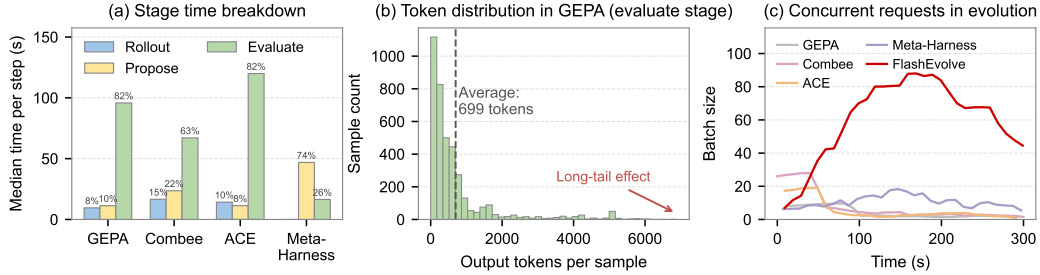


Figure 2: Profiling results of inefficiency in synchronized agent evolution. (a) Stage execution is serial, and stage time is highly imbalanced. (b) Within a single stage, output lengths show a long-tail distribution, so the slowest requests determine stage completion time. (c) Serial stage execution and intra-stage imbalance reduce effective concurrency, while FlashEvolve keeps more requests in flight.

72 this beyond the agent itself, evolving the generated programs the agent uses to solve problems, where  
 73 the LLM acts as a mutation operator and an external evaluator scores each candidate.

74 An agent evolve loop iterates over multiple iteration steps, where each step consists of several stages,  
 75 as illustrated in Figure 1. The LLM-heavy stages are typically *Generate*, *Propose*, and *Evaluate*. The  
 76 *Generate* stage runs the current artifact on tasks to collect trajectories. The *Propose* stage reflects on  
 77 these trajectories to produce a new candidate artifact. The *Evaluate* stage scores the candidate against  
 78 task signals and filters out underperforming ones. A subsequent update commits the new artifact to  
 79 the artifact pool. At the start of each step, new candidate artifacts are selected from the pool, through  
 80 methods like Pareto-aware sampling [2] or evolutionary tournaments [20].

## 81 2.2 Inefficiency in Agent Evolution: Sequential and Imbalanced Stages

82 Despite its algorithmic appeal, agent evolution remains expensive in wall-clock time. Based on  
 83 our experiments, even with state-of-the-art LLM serving infrastructure such as vLLM [12], which  
 84 supports continuous batching and prefix caching, GEPA with Qwen3-8B takes 50 minutes to complete  
 85 49 evolution steps on IFBench [23], and 134 minutes to complete 411 steps on HotpotQA [31].

86 This inefficiency stems from sequential and synchronized stage execution. Each evolution step runs  
 87 its LLM-heavy stages serially, and each stage internally waits for all parallel LLM requests to finish  
 88 before advancing to the next stage. This structure produces two compounding costs. First, the serial  
 89 chain forces total step time to be the sum of per-stage durations, with no opportunity to overlap  
 90 stages. As shown in Figure 2(a), stage time is highly imbalanced, so different stages can become  
 91 the bottleneck depending on the workload and algorithm. Second, the synchronization barrier at  
 92 each stage’s end forces the entire batch to wait for the slowest one. As shown in Figure 2(b), output  
 93 lengths within a stage show a long-tail distribution, so a small number of long requests determine  
 94 stage completion time. Consequently, sequential execution and intra-stage imbalance reduce effective  
 95 concurrency and leave the LLM backend underutilized, as shown in Figure 2(c).

96 Such inefficiency cannot be solved by simply launching more LLM requests in parallel. Agent  
 97 evolution must convert a synchronized multi-stage loop into a streaming workflow while preserving  
 98 artifact-evolution semantics. This creates two challenges. First, asynchrony introduces artifact-level  
 99 staleness: intermediate results may be produced from an artifact pool that has already changed before  
 100 they are consumed. Second, naive parallel scaling can amplify workload imbalance: fast stages may  
 101 overproduce items for slow stages, while long-tail requests within a stage can still delay downstream  
 102 execution. This causes queue buildup, longer staleness windows, and wasted LLM work. These  
 103 challenges require orchestration mechanisms that jointly manage staleness and workload balance.

104 **Analogy to Asynchronous RL.** These challenges are related to synchronous LLM RL sys-  
 105 tems [25, 1, 8], which also suffer from synchronization overhead and workload imbalance. Asyn-  
 106 chronous RL addresses this by overlapping rollout generation with training and controlling off-policy  
 107 optimization [5, 38, 24]. Agent evolution differs in two key ways. First, it contains multiple LLM  
 108 inference stages rather than a single "rollout" stage in RL. Each stage has batched generation behavior  
 109 and its own long-tail imbalance. Second, staleness occurs over inspectable language artifacts, such as

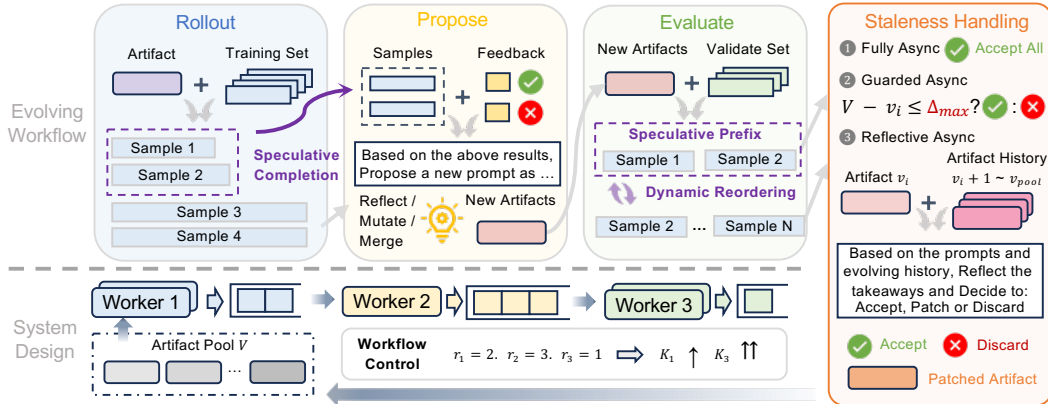


Figure 3: Overview of FlashEvolve. FlashEvolve executes agent evolution with asynchronous workers and queues across stages. Workers pass partial or completed results through queues instead of waiting for a whole stage to finish. FlashEvolve further uses speculative completion, validation-set reordering, workflow control, and staleness-aware handling to improve throughput while limiting data staleness.

110 prompts, memories, harness code, and programs, rather than continuous model weights. This allows  
 111 a more flexible design space for staleness handling policies.

### 112 3 FlashEvolve: Asynchronous Framework for Agent Evolution

113 We present FlashEvolve, an asynchronous framework that removes the sequential and imbalanced  
 114 behavior identified in Section 2.2. FlashEvolve decomposes an evolution loop into asynchronous  
 115 workers connected by queues, so different stages and evolution steps can overlap. Each queue item  
 116 carries the artifact state and pool version, allowing FlashEvolve to detect stale items. On top of this  
 117 execution model, FlashEvolve introduces staleness-aware data handling, speculative stage completion,  
 118 and adaptive workflow control to improve the time efficiency of evolution.

#### 119 3.1 Asynchronous Execution with Workers and Queues

120 **Asynchronous workers.** FlashEvolve turns a synchronized evolution step into asynchronous workers  
 121 connected by queues. Instead of waiting for artifact proposal, validation, and pool update to finish  
 122 before starting the next step, workers continuously process ready items and pass their outputs to  
 123 downstream queues. This allows different stages and evolution steps to overlap. Each stage has an  
 124 input queue and a set of workers. A queue item carries the artifact being evolved, the input/output,  
 125 and the artifact-pool version  $v_i$  at item creation. The pool version increases after each pool update, so  
 126 FlashEvolve can compare  $v_i$  with the current version  $v$  to detect stale items.

127 **Worker concurrency.** To improve system throughput, FlashEvolve assigns a worker count  $K_i$  to  
 128 each asynchronous stage  $i$ . A larger  $K_i$  allows more tasks in stage  $i$  to issue LLM requests at the  
 129 same time, which increases per-stage concurrency so the whole pipeline is not bottlenecked by the  
 130 throughput of a single slow or imbalanced stage. The tradeoff is data staleness: larger worker counts  
 131 increase the chance that queued items were generated from an older artifact pool state.

#### 132 3.2 Staleness-Aware Data Handling

133 FlashEvolve supports three policies for handling such stale items with different tradeoffs:

- 134 • *Full Async* does not check artifact pool versions and allows all items to continue through the  
 135 pipeline. This policy preserves all completed work and maximizes throughput, but stale items may  
 136 introduce outdated updates into the artifact pool and impact convergence.
- 137 • *Guarded Async* discards an item when its version gap exceeds a threshold  $\Delta_{\max}$ . Let  $v_i$  denote  
 138 the artifact-pool version used to generate item  $i$ , and let  $v$  denote the current artifact-pool version.  
 139 The version gap is defined as  $\Delta_i = v - v_i$ . Guarded Async allows item  $i$  to continue only when

140  $\Delta_i \leq \Delta_{\max}$ ; otherwise, it discards the item. This policy prevents highly stale items, but will waste  
141 the generated tokens that already spent on discarded items.

142 • *Reflective Async* inspects and updates stale items by adding a new reflection worker stage. For  
143 an item  $i$  with version gap  $\Delta_i > 0$ , the reflection worker uses the stale item and all artifact-pool  
144 updates between versions  $v_i$  and  $v$  to decide whether the item still contributes a useful change. If so,  
145 it patches the item against the current artifact pool state and lets it continue; otherwise, FlashEvolve  
146 discards it. Non-stale items continue without reflection. This policy avoids uncontrolled stale  
147 updates while reusing useful stale items, reducing wasted LLM generations.

148 **Why language-space staleness can be repaired.** Language-space staleness is discrete and in-  
149 spectable, unlike parameter staleness in asynchronous RL, which is continuous and opaque. In RL, a  
150 stale item is tied to an older point in weight space, so systems typically handle it through importance  
151 weighting, bounded delay, or discard. In agent evolution, a stale item is text or code, such as a  
152 prompt edit, memory update, harness mutation, or generated program. FlashEvolve can therefore  
153 inspect the stale item together with the intervening artifact history and decide whether the edit is  
154 orthogonal, already subsumed, or conflicting with the current artifact pool. This makes repair a  
155 first-class operation: stale items can be patched when they contain reusable information, or discarded  
156 when they are too specific or inconsistent. Figure 5 shows an example where FlashEvolve filters  
157 task-specific stale content and keeps transferable principles to form a compact prompt patch.

### 158 3.3 Speculative Stage Completion

159 Asynchronous workers remove waiting between stages, but each worker may still wait for all LLM  
160 requests in its current stage before writing to the next queue. This still creates an intra-stage  
161 synchronization barrier, especially in rollout and evaluate stages where a minibatch contains many  
162 LLM requests. To reduce this barrier, FlashEvolve allows a stage to release partial output after a  
163 fraction  $\alpha_{\text{spec}}^i \in (0, 1]$  of its requests has finished. The worker packages the completed results as  
164 a tentative queue item and continues the remaining requests in the background, while downstream  
165 workers can start from the tentative item.

166 For rollout, this means completed samples can be forwarded as soon as they are available. For  
167 evaluation, FlashEvolve adds a score threshold to avoid forwarding weak candidates. After the first  
168  $\alpha_{\text{spec}}^i$  fraction of evaluation requests finishes, the worker computes a partial score. If the partial score  
169 exceeds the current pool score, FlashEvolve inserts the candidate into the pool as a *speculative artifact*.  
170 When full evaluation finishes, the artifact is confirmed if it still passes the acceptance condition;  
171 otherwise, it is removed. If a speculative artifact is later removed, downstream items derived from it  
172 are marked stale and handled by the same staleness-aware policy in Section 3.2; they cannot update  
173 the confirmed pool without passing the normal validation path.

174 **Validation-set reordering.** Speculative completion is more reliable when the early validation samples  
175 are informative. We call the first  $\alpha_{\text{spec}}$  fraction of the validation set the *speculative prefix*. FlashEvolve  
176 reorders the validation set using sample pass history: samples that pass for  $w$  consecutive rounds are  
177 moved out of the speculative prefix and placed later in the validation order. This keeps easy samples  
178 from dominating the early signal and leaves more discriminative samples in the prefix. We set  $w = 3$   
179 to avoid reacting to one-round noise while keeping the prefix responsive to artifact improvement.

### 180 3.4 Adaptive Workflow Control

181 Different stages in an evolution loop produce and consume items at different rates. A stage with  
182 short LLM requests can quickly fill the queue of a later stage whose requests are longer or more  
183 imbalanced. If workers keep running at a fixed concurrency, the queue keeps growing and many  
184 items become stale before they are processed. FlashEvolve therefore monitors queue pressure and  
185 version gap to adjust worker behavior, making execution more balanced and token efficient.

186 **Adaptive worker reallocation.** FlashEvolve measures the item production rate of each asynchronous  
187 stage. The production rate is the number of queue items that a stage writes to its downstream queue  
188 per second. A stage with a much lower production rate can limit the whole workflow, while a stage  
189 with a much higher production rate can overfeed downstream queues.

190 FlashEvolve compares production rates across stages and adjusts their worker counts. If a stage  
191 produces items at less than half the median stage rate, we increase its worker count. If a stage produces

Table 1: Throughput comparison on GEPA workloads. LLM throughput measures the output token rate of the whole system. Proposal throughput measures the rate of new candidate artifact generation.

Method	LLM Throughput (token/s)				Proposal Throughput (proposal/min)			
	IFBench	HotpotQA	HoVer	AIME	IFBench	HotpotQA	HoVer	AIME
<i>vLLM with Qwen3-8B</i>								
GEPA	963	30	461	200	1.9	4.6	2.5	2.2
Combee ( $K=10$ )	696	38	810	994	1.2	2.7	2.0	6.2
Combee ( $K=40$ )	900	44	891	977	0.7	4.5	2.0	1.6
FlashEvolve	<b>2,688</b>	<b>93</b>	<b>1,255</b>	<b>998</b>	<b>8.9</b>	<b>8.8</b>	<b>5.9</b>	<b>11.4</b>
<i>API with GPT-4o-mini</i>								
GEPA	361	14	142	103	1.7	2.4	1.8	1.3
Combee ( $K=10$ )	397	18	348	211	1.0	1.4	0.8	1.0
Combee ( $K=40$ )	389	23	214	336	0.8	1.2	0.7	0.6
FlashEvolve	<b>791</b>	<b>32</b>	<b>352</b>	<b>485</b>	<b>10.1</b>	<b>8.0</b>	<b>9.1</b>	<b>6.6</b>

192 items at more than twice the median stage rate, we decrease its worker count. Each adjustment  
 193 changes the worker count by at most one, and each stage has a minimum and maximum worker count.  
 194 This avoids large swings while still correcting persistent throughput imbalance.

### 195 3.5 Implementation

196 FlashEvolve is implemented in Python with lightweight threads and in-process queues. Each stage  
 197 is executed by a small worker pool, queue items carry the artifact state and pool version, and pool  
 198 updates are applied under a lock. For a fair comparison, we run all open-source baselines and  
 199 FlashEvolve on the same LLM serving stack: the native LLM calls in different algorithms are  
 200 replaced by the same DSPy client backed by a local vLLM [12] server with an OpenAI-compatible  
 201 endpoint. Thus all methods benefit from the same continuous batching and KV-cache reuse, and  
 202 throughput differences mainly reflect the optimization of the evolution pipeline. The same interface  
 203 is also used for API-based experiments by changing only the endpoint and model name.

## 204 4 Evaluation

### 205 4.1 Experimental Setup

206 **Evolving algorithm baselines.** We evaluate FlashEvolve on three test-time evolution algorithms that  
 207 optimize different artifacts. We use GEPA [2], which evolves prompts through execution feedback  
 208 and reflection, as the main algorithm for end-to-end comparison and ablation studies. We set the  
 209 rollout minibatch size  $mb = 3$ , following the default setting. We also reproduce Combee [15] as a  
 210 scaling-oriented baseline. Combee scales batch-level parallelism to improve throughput. We evaluate  
 211 its reported fixed-batch variants with parallelization sizes  $B = 10$  and  $40$ .

212 We also evaluate FlashEvolve on ACE [34] and Meta-Harness [14]. ACE evolves agent context  
 213 playbooks, while Meta-Harness evolves harness code. These algorithms use different artifact types,  
 214 but all fall into the abstraction optimized by FlashEvolve: a multi-stage evolution loop with LLM-  
 215 heavy stages, queueable intermediate results, and a shared artifact pool that is updated over time.

216 **Models and deployment.** For open-source model experiments, we use Qwen3-8B [30], which is the  
 217 default model used in GEPA and provides a representative setting for studying test-time evolution  
 218 behavior. We serve Qwen3-8B with vLLM on a single NVIDIA H100 80GB GPU and an AMD EPYC  
 219 9534 CPU. For API-based experiments, we use GPT-4o-mini [9], which shows similar evolution  
 220 behavior to the open-source setting while representing a common remote-serving deployment.

221 **Datasets.** We use the benchmark datasets used in each original evolution algorithm. For GEPA,  
 222 we evaluate on IFBench [23] for instruction following, HotpotQA [31] for knowledge retrieval,  
 223 HoVer [11] for multi-hop verification, and AIME for mathematical reasoning. For ACE, we use  
 224 FiNER [17] and FormulaReasoning, which are finance and numerical-reasoning datasets, respectively.  
 225 For Meta-Harness, we use a mixture of Symptom2Disease for medical diagnosis classification and  
 226 AGNews [35] for topic categorization. These datasets cover diverse domains of agent applications.

Table 2: Validation score(%) and normalized evolution rate within 30 minutes on GEPA workloads using Qwen3-8B. Evolution rate measures the score improvement achieved within the time budget, normalized by the improvement achieved by serial GEPA, reflecting the speedup of evolution progress. AIME reports “-” because serial GEPA shows no score improvement within the budget.

Method	IFBench		HoVer		HotpotQA		AIME	
	Score	Norm. Rate	Score	Norm. Rate	Score	Norm. Rate	Score	Norm. Rate
GEPA	87.6	1.00	39.8	1.00	<b>63.3</b>	<b>1.00</b>	10.0	-
Combee ( $K=10$ )	88.5	1.39	41.2	1.09	62.5	0.94	10.0	-
Combee ( $K=40$ )	86.5	0.55	40.5	1.05	58.6	0.63	10.0	-
FlashEvolve	<b>90.6</b>	<b>2.27</b>	<b>42.0</b>	<b>1.15</b>	61.7	0.88	<b>15.0</b>	-

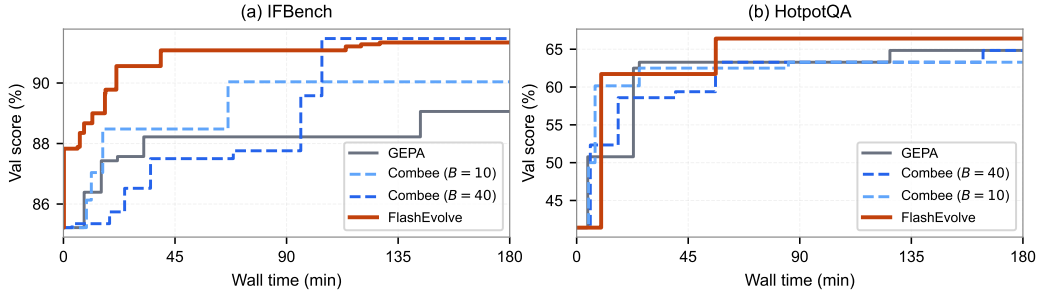


Figure 4: Longer-time validation score evolution over wall-clock time with Qwen3-8B.

## 227 4.2 Improvement on GEPA

228 **System throughput improvement.** Table 1 first shows that FlashEvolve substantially improves  
 229 LLM throughput. This indicates that asynchronous stage orchestration keeps the LLM backend  
 230 busier by overlapping requests from different stages and evolution steps. On local vLLM serving,  
 231 FlashEvolve improves LLM throughput by  $3.4\times$  on average over GEPA and  $1.9\times$  on average over  
 232 the best Combee setting. The same trend holds for API-based serving: FlashEvolve improves LLM  
 233 throughput by  $2.9\times$  on average over GEPA and  $1.5\times$  on average over the best Combee setting.

234 We also show that higher LLM throughput translates into faster artifact exploration. On local vLLM  
 235 serving, FlashEvolve improves proposal throughput by  $3.5\times$  on average over GEPA and  $3.5\times$  on  
 236 average over the best Combee setting. On API-based serving, FlashEvolve improves proposal  
 237 throughput by  $4.9\times$  on average over GEPA and  $8.4\times$  on average over the best Combee setting.  
 238 Across all settings, FlashEvolve sustains more than 5.9 proposals/min and up to 11.4 proposals/min,  
 239 showing that it substantially increases the rate at which evolution tests new candidate artifacts.

240 **Evolving efficiency improvement.** Table 2 reports validation score and normalized evolution rate  
 241 within a fixed 30-minute budget. Across the three workloads where the GEPA baseline makes  
 242 measurable progress, FlashEvolve achieves an average normalized evolution rate of  $1.43\times$ . The  
 243 strongest gain appears on IFBench, where FlashEvolve improves the validation score from 87.6 to  
 244 90.6 and reaches a  $2.27\times$  normalized evolution rate. On HoVer, FlashEvolve also achieves the best  
 245 score and a  $1.15\times$  normalized rate. On HotpotQA, FlashEvolve does not reach the best validation  
 246 score within the 30-minute budget, but Figure 4 shows that this advantage emerges under a longer  
 247 budget. On AIME, both GEPA and Combee remain at the initial score of 10.0%, while FlashEvolve  
 248 reaches 15.0%, making it the only method that improves over the initial score.

249 **Long-time evolution.** Figure 4 reports validation score over a longer 180-minute budget. FlashEvolve  
 250 reaches strong validation scores earlier than the synchronous baselines. On IFBench, FlashEvolve  
 251 reaches 91% in 39.3 minutes, while Combee ( $B=40$ ) reaches the same score region after 104.2  
 252 minutes and eventually approaches a similar final score. On HotpotQA, FlashEvolve reaches its best  
 253 score of 66.41% at 56.1 minutes and maintains the highest validation score over the full budget, while  
 254 all baselines remain below 65%. These results show that asynchronous evolution accelerates useful  
 255 artifact discovery, and in some workloads also improves the final score under a longer time budget.

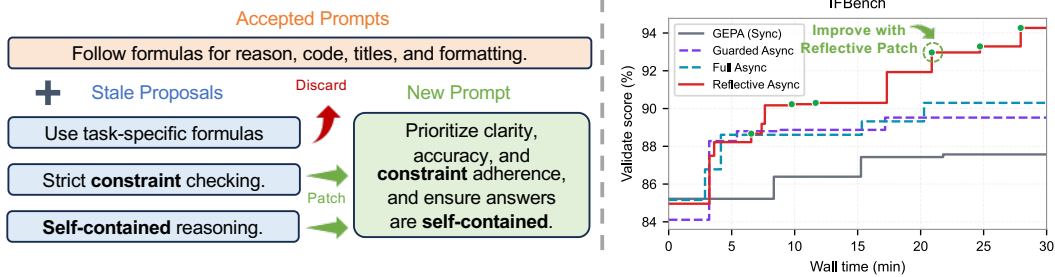


Figure 5: Staleness handling on IFBench with Qwen3-8B. The left panel shows an example of reflective prompt repair. The right panel compares Full Async, Guarded Async, and Reflective Async.

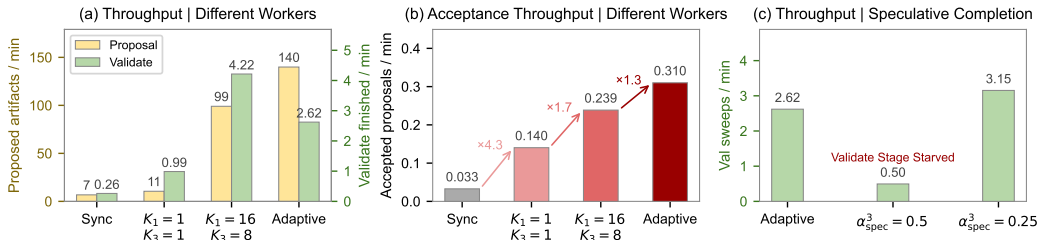


Figure 6: Ablation of worker concurrency and speculative completion on IFBench with Qwen3-8B. (a) Worker allocation varies throughput. (b) Adaptive worker control achieves the highest accepted proposal throughput by balancing proposal generation and validation. (c) Speculative completion improves validation throughput when the prefix threshold is properly set.

### 256 4.3 Ablation Studies

257 **Comparison across staleness handling methods.** Figure 5 compares three staleness handling  
 258 policies on IFBench. Full Async and Guarded Async achieve similar final scores, but their behaviors  
 259 differ. Full Async preserves all stale items and therefore keeps high throughput, while Guarded Async  
 260 discards highly stale items to avoid outdated updates. In this case, the two methods perform similarly.

261 Reflective Async achieves the best evolution efficiency. It reaches a validation score of 94.3% within  
 262 the 30-minute budget. This shows that stale items are not always useless: when the stale artifact is  
 263 text, FlashEvolve can inspect it, discard task-specific noise, and reuse transferable principles. Figure 5  
 264 also shows a simplified repair example extracted from our logs. FlashEvolve discards task-specific  
 265 formulas because the accepted prompt already contains general instruction-following rules and the  
 266 formula does not transfer across tasks. In contrast, new takeaways such as stricter constraint checking  
 267 and self-contained reasoning are distilled into a compact prompt patch. We also observe that many  
 268 score jumps in the Reflective Async curve come from prompts after such patches, suggesting that  
 269 reflective repair improves the quality of prompt proposals rather than only increasing throughput.

270 **Concurrent Workers.** We study how worker counts affect stage throughput. As shown in Figure 6(a),  
 271 larger worker counts greatly increase proposal throughput, from 7 artifacts/min in the synchronous  
 272 setting to 99 artifacts/min with  $K_1=16, K_3=8$ . However, validation throughput does not scale  
 273 uniformly, showing that naive scaling can shift the bottleneck across stages. Adaptive control  
 274 balances queue pressure and stage rates; although it reduces validation throughput compared with  
 275 the fixed large-worker setting, it achieves the highest accepted proposal throughput in Figure 6(b),  
 276 suggesting that a more balanced worker allocation can produce a higher rate of high-quality candidates  
 277 rather than merely increasing raw proposal volume.

278 **Speculative Stage Completion.** Figure 6(c) shows that speculative completion can improve validation  
 279 throughput when the prefix threshold is properly set. With  $\alpha_{spec}^3=0.25$ , validation-stage throughput  
 280 increases to 3.15 validations/min and the validation score improves by 4.49 percentage points within  
 281 30 minutes. However, when  $\alpha_{spec}^3=0.5$ , the speculative gate becomes less effective: candidates must  
 282 wait for a larger partial validation prefix before being released, which reduces early handoff and

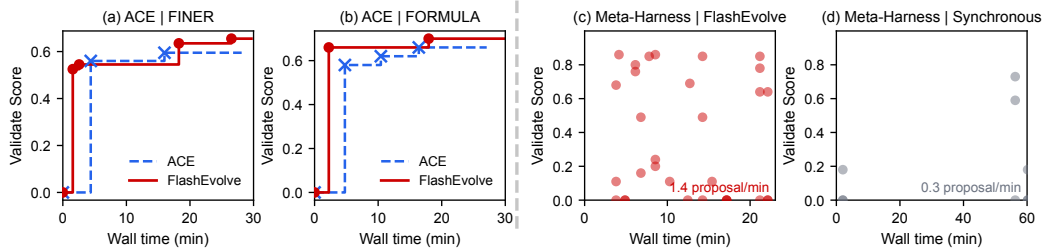


Figure 7: FlashEvolve on other algorithms for agent evolution. (a) and (b) compare validation score evolution curve on FiNER and Formula over a fixed time budget. (c) and (d) compare Meta-Harness proposal rate and validation scores of different proposals on Symptom2Disease and AgNews.

lowers validation throughput. Overall, speculative completion can be useful in some settings, but its accuracy and efficiency depend on  $\alpha_{\text{spec}}$  and dataset characteristics. We therefore treat it as an optional optimization in FlashEvolve rather than include it in the main evaluation.

#### 4.4 Improvement on Other Algorithms: ACE and Meta-Harness

FlashEvolve is algorithm-agnostic. It does not rely on a specific artifact type, but only assumes that the evolution loop contains multiple stages that need orchestration. We evaluate FlashEvolve on two other algorithms including ACE [34] and Meta-Harness [14], which evolves harness code.

Figure 7(a)(b) shows the wall-clock evolution curves of ACE. FlashEvolve reaches better validation scores within the same 30-minute budget on both tasks. The validation score improves from 0.6 to 0.66 on FiNER and from 0.66 to 0.7 on Formula, demonstrating higher efficiency.

Figure 7(c)(d) compares FlashEvolve with the synchronous Meta-Harness baseline. FlashEvolve improves the proposal and validation throughput from 0.3 to 1.4 proposals/min, giving a  $4.7\times$  speedup. Since the open-source model has relatively weak code-generation capability, harness-code evolution progresses slowly in both settings. We therefore report the score distribution of different proposed harnesses. With higher proposal throughput, FlashEvolve samples and validates more harness candidates within the same time budget, leading to a higher potential of improvement.

## 5 Conclusion

This paper presents FlashEvolve, a framework for accelerating agent evolution in wall-clock time. FlashEvolve replaces synchronous stage execution with asynchronous workers and queues, allowing LLM-heavy stages and evolution steps to overlap. It preserves evolution semantics through artifact-pool versioning and staleness-aware policies that update, discard, or patch stale artifacts, and further improves efficiency with speculative stage completion and adaptive workflow control. On GEPA workloads, FlashEvolve achieves  $3.5\times$  higher proposal throughput over the synchronous implementation on local vLLM serving. The same execution model also generalizes to context evolution with ACE and harness-code evolution with Meta-Harness,

**Limitations.** FlashEvolve currently supports only a limited set of agent-evolution algorithms, and each integration still requires algorithm-specific implementation effort. Although the worker and queue abstraction is general, mapping a new algorithm to this abstraction requires implementing its stages, queue items, artifact state, and update rules. Our current evaluation also focuses on representative prompt, context, and harness-code evolution workloads, and broader coverage of evolution algorithms and artifact types remains future work.

**Future Work.** In future work, we plan to expand FlashEvolve with a more general plugin interface for defining stages, artifacts, staleness policies, and evaluation logic, so that new evolution algorithms can be integrated with less manual engineering. We also plan to extend FlashEvolve to more types of artifact evolution, such as memory, tool-use policies and generated programs.

318 **References**

- 319 [1] Nemo rl: A scalable and efficient post-training library. [https://github.com/NVIDIA-NeMo/](https://github.com/NVIDIA-NeMo/RL)  
320 RL, 2025. GitHub repository.
- 321 [2] L. A. Agrawal, S. Tan, D. Soylu, N. Ziemis, R. Khare, K. Opsahl-Ong, A. Singhvi, H. Shandilya,  
322 M. J. Ryan, M. Jiang, et al. Gepa: Reflective prompt evolution can outperform reinforcement  
323 learning. *arXiv preprint arXiv:2507.19457*, 2025.
- 324 [3] H. Assumpção, D. Ferreira, L. Campos, and F. Murai. Codeevolve: An open source evolutionary  
325 coding agent for algorithm discovery and optimization. *arXiv preprint arXiv:2510.14150*, 2025.
- 326 [4] J. Fang, Y. Peng, X. Zhang, Y. Wang, X. Yi, G. Zhang, Y. Xu, B. Wu, S. Liu, Z. Li, et al. A  
327 comprehensive survey of self-evolving ai agents: A new paradigm bridging foundation models  
328 and lifelong agentic systems. *arXiv preprint arXiv:2508.07407*, 2025.
- 329 [5] W. Fu, J. Gao, X. Shen, C. Zhu, Z. Mei, C. He, S. Xu, G. Wei, J. Mei, J. Wang, et al. Areal: A  
330 large-scale asynchronous reinforcement learning system for language reasoning. *arXiv preprint*  
331 *arXiv:2505.24298*, 2025.
- 332 [6] H.-a. Gao, J. Geng, W. Hua, M. Hu, X. Juan, H. Liu, S. Liu, J. Qiu, X. Qi, Y. Wu, et al. A survey  
333 of self-evolving agents: On path to artificial super intelligence. *arXiv preprint arXiv:2507.21046*,  
334 1, 2025.
- 335 [7] D. Guo, D. Yang, H. Zhang, J. Song, P. Wang, Q. Zhu, R. Xu, R. Zhang, S. Ma, X. Bi, et al.  
336 Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv*  
337 *preprint arXiv:2501.12948*, 2025.
- 338 [8] Z. Hu, H. Ouyang, C. Chen, Z. Pan, Y. Guan, Z. Yu, Z. Wang, S. Swanson, and Y. Ding.  
339 Jigsawrl: Assembling rl pipelines for efficient llm post-training, 2026. URL [https://arxiv.](https://arxiv.org/abs/2604.23838)  
340 [org/abs/2604.23838](https://arxiv.org/abs/2604.23838).
- 341 [9] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda,  
342 A. Hayes, A. Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- 343 [10] A. Jaech, A. Kalai, A. Lerer, A. Richardson, A. El-Kishky, A. Low, A. Helyar, A. Madry,  
344 A. Beutel, A. Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- 345 [11] Y. Jiang, S. Bordia, Z. Zhong, C. Dognin, M. Singh, and M. Bansal. Hover: A dataset for many-  
346 hop fact extraction and claim verification. In *Findings of the Association for Computational*  
347 *Linguistics: EMNLP 2020*, pages 3441–3460, 2020.
- 348 [12] W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. Gonzalez, H. Zhang, and I. Stoica.  
349 Efficient memory management for large language model serving with pagedattention. In  
350 *Proceedings of the 29th symposium on operating systems principles*, pages 611–626, 2023.
- 351 [13] R. T. Lange, Y. Imajuku, and E. Cetin. Shinkaevolve: Towards open-ended and sample-efficient  
352 program evolution. *arXiv preprint arXiv:2509.19349*, 2025.
- 353 [14] Y. Lee, R. Nair, Q. Zhang, K. Lee, O. Khattab, and C. Finn. Meta-harness: End-to-end  
354 optimization of model harnesses. *arXiv preprint arXiv:2603.28052*, 2026.
- 355 [15] H. Li, R. He, Q. Zhang, C. Ji, Q. Mang, X. Chen, L. A. Agrawal, W.-L. Liao, E. Yang, A. Cheung,  
356 et al. Combee: Scaling prompt learning for self-improving language model agents. *arXiv*  
357 *preprint arXiv:2604.04247*, 2026.
- 358 [16] X. Lou, M. Lázaro-Gredilla, A. Dedieu, C. Wendelken, W. Lehrach, and K. P. Murphy. Auto-  
359 harness: improving llm agents by automatically synthesizing a code harness. *arXiv preprint*  
360 *arXiv:2603.03329*, 2026.
- 361 [17] L. Loukas, M. Fergadiotis, I. Chalkidis, E. Spyropoulou, P. Malakasiotis, I. Androutsopoulos,  
362 and G. Paliouras. Finer: Financial numeric entity recognition for xbrl tagging. In *Proceedings*  
363 *of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
364 *Papers)*, pages 4419–4431, 2022.

- 365 [18] H. Lu, H. Huang, Y. Zhou, C. Li, and N. Zhu. Empirical-mcts: Continuous agent evolution via  
366 dual-experience monte carlo tree search. *arXiv preprint arXiv:2602.04248*, 2026.
- 367 [19] A. Madaan, N. Tandon, P. Gupta, S. Hallinan, L. Gao, S. Wiegrefe, U. Alon, N. Dziri,  
368 S. Prabhunoye, Y. Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in*  
369 *neural information processing systems*, 36:46534–46594, 2023.
- 370 [20] A. Novikov, N. Vū, M. Eisenberger, E. Dupont, P.-S. Huang, A. Z. Wagner, S. Shirobokov,  
371 B. Kozlovskii, F. J. Ruiz, A. Mehrabian, et al. Alphaevolve: A coding agent for scientific and  
372 algorithmic discovery. *arXiv preprint arXiv:2506.13131*, 2025.
- 373 [21] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal,  
374 K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback.  
375 *Advances in neural information processing systems*, 35:27730–27744, 2022.
- 376 [22] S. Ouyang, J. Yan, I. Hsu, Y. Chen, K. Jiang, Z. Wang, R. Han, L. T. Le, S. Daruki, X. Tang,  
377 et al. Reasoningbank: Scaling agent self-evolving with reasoning memory. *arXiv preprint*  
378 *arXiv:2509.25140*, 2025.
- 379 [23] V. Pyatkin, S. Malik, V. Graf, H. Ivison, S. Huang, P. Dasigi, N. Lambert, and H. Hajishirzi.  
380 Generalizing verifiable instruction following. *arXiv preprint arXiv:2507.02833*, 2025.
- 381 [24] G. Sheng, Y. Tong, B. Wan, W. Zhang, C. Jia, X. Wu, Y. Wu, X. Li, C. Zhang, Y. Peng, et al.  
382 Laminar: A scalable asynchronous rl post-training framework. *arXiv preprint arXiv:2510.12633*,  
383 2025.
- 384 [25] G. Sheng, C. Zhang, Z. Ye, X. Wu, W. Zhang, R. Zhang, Y. Peng, H. Lin, and C. Wu. Hybridflow:  
385 A flexible and efficient rlhf framework. In *Proceedings of the Twentieth European Conference*  
386 *on Computer Systems*, pages 1279–1297, 2025.
- 387 [26] N. Shinn, F. Cassano, A. Gopinath, K. Narasimhan, and S. Yao. Reflexion: Language agents  
388 with verbal reinforcement learning. *Advances in neural information processing systems*, 36:  
389 8634–8652, 2023.
- 390 [27] M. Shoeybi, M. Patwary, R. Puri, P. LeGresley, J. Casper, and B. Catanzaro. Megatron-lm:  
391 Training multi-billion parameter language models using model parallelism. *arXiv preprint*  
392 *arXiv:1909.08053*, 2019.
- 393 [28] X. Wang, C. Li, Z. Wang, F. Bai, H. Luo, J. Zhang, N. Jojic, E. P. Xing, and Z. Hu. Promptagent:  
394 Strategic planning with language models enables expert-level prompt optimization. *arXiv*  
395 *preprint arXiv:2310.16427*, 2023.
- 396 [29] E. Xiao, Y. Zeng, A. Chen, C.-J. Li, A. Bertsch, and G. Neubig. Prompt-mii: Meta-learning  
397 instruction induction for llms. *arXiv preprint arXiv:2510.16932*, 2025.
- 398 [30] A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv, et al.  
399 Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- 400 [31] Z. Yang, P. Qi, S. Zhang, Y. Bengio, W. Cohen, R. Salakhutdinov, and C. D. Manning. Hotpotqa:  
401 A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018*  
402 *conference on empirical methods in natural language processing*, pages 2369–2380, 2018.
- 403 [32] M. Yuksekgonul, F. Bianchi, J. Boen, S. Liu, Z. Huang, C. Guestrin, and J. Zou. Textgrad:  
404 "Automatic" differentiation" via text. *arXiv preprint arXiv:2406.07496*, 2024.
- 405 [33] G. Zhang, H. Ren, C. Zhan, Z. Zhou, J. Wang, H. Zhu, W. Zhou, and S. Yan. Memevolve:  
406 Meta-evolution of agent memory systems. *arXiv preprint arXiv:2512.18746*, 2025.
- 407 [34] Q. Zhang, C. Hu, S. Upasani, B. Ma, F. Hong, V. Kamanuru, J. Rainton, C. Wu, M. Ji, H. Li,  
408 et al. Agentic context engineering: Evolving contexts for self-improving language models.  
409 *arXiv preprint arXiv:2510.04618*, 2025.
- 410 [35] X. Zhang, J. Zhao, and Y. LeCun. Character-level convolutional networks for text classification.  
411 *Advances in neural information processing systems*, 28, 2015.

- 412 [36] Y. Zhao, A. Gu, R. Varma, L. Luo, C.-C. Huang, M. Xu, L. Wright, H. Shojanazeri, M. Ott,  
413 S. Shleifer, et al. Pytorch fsdp: experiences on scaling fully sharded data parallel. *arXiv preprint*  
414 *arXiv:2304.11277*, 2023.
- 415 [37] L. Zheng, L. Yin, Z. Xie, C. Sun, J. Huang, C. H. Yu, S. Cao, C. Kozyrakis, I. Stoica, J. E.  
416 Gonzalez, et al. Sglang: Efficient execution of structured language model programs. *Advances*  
417 *in neural information processing systems*, 37:62557–62583, 2024.
- 418 [38] Y. Zhong, Z. Zhang, X. Song, H. Hu, C. Jin, B. Wu, N. Chen, Y. Chen, Y. Zhou, C. Wan, et al.  
419 Streamrl: Scalable, heterogeneous, and elastic rl for llms with disaggregated stream generation.  
420 *arXiv preprint arXiv:2504.15930*, 2025.

## 421 **NeurIPS Paper Checklist**

### 422 **1. Claims**

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

425 Answer: [Yes]

426 Justification: The abstract and introduction state that FlashEvolve improves wall-clock  
427 efficiency for agent evolution through asynchronous stage orchestration, staleness-aware  
428 handling, speculative completion, and adaptive workflow control. These claims are supported  
429 by the experimental results on GEPA, Combee, ACE, and Meta-Harness.

430 Guidelines:

- 431 • The answer [N/A] means that the abstract and introduction do not include the claims  
432 made in the paper.
- 433 • The abstract and/or introduction should clearly state the claims made, including the  
434 contributions made in the paper and important assumptions and limitations. A [No] or  
435 [N/A] answer to this question will not be perceived well by the reviewers.
- 436 • The claims made should match theoretical and experimental results, and reflect how  
437 much the results can be expected to generalize to other settings.
- 438 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
439 are not attained by the paper.

### 440 **2. Limitations**

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

442 Answer: [Yes]

443 Justification: The paper includes a limitations discussion after the conclusion. It discusses  
444 the current scope of supported evolution algorithms and the dependence on customized  
445 implementations for each target codebase.

446 Guidelines:

- 447 • The answer [N/A] means that the paper has no limitation while the answer [No] means  
448 that the paper has limitations, but those are not discussed in the paper.
- 449 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 450 • The paper should point out any strong assumptions and how robust the results are to  
451 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
452 model well-specification, asymptotic approximations only holding locally). The authors  
453 should reflect on how these assumptions might be violated in practice and what the  
454 implications would be.
- 455 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
456 only tested on a few datasets or with a few runs. In general, empirical results often  
457 depend on implicit assumptions, which should be articulated.
- 458 • The authors should reflect on the factors that influence the performance of the approach.  
459 For example, a facial recognition algorithm may perform poorly when image resolution  
460 is low or images are taken in low lighting. Or a speech-to-text system might not be  
461 used reliably to provide closed captions for online lectures because it fails to handle  
462 technical jargon.
- 463 • The authors should discuss the computational efficiency of the proposed algorithms  
464 and how they scale with dataset size.
- 465 • If applicable, the authors should discuss possible limitations of their approach to  
466 address problems of privacy and fairness.
- 467 • While the authors might fear that complete honesty about limitations might be used by  
468 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
469 limitations that aren’t acknowledged in the paper. The authors should use their best  
470 judgment and recognize that individual actions in favor of transparency play an impor-  
471 tant role in developing norms that preserve the integrity of the community. Reviewers  
472 will be specifically instructed to not penalize honesty concerning limitations.

### 473 **3. Theory assumptions and proofs**

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

476 Answer: [N/A]

477 Justification: The paper does not present formal theoretical results, theorems, or proofs. Its  
478 contributions are systems mechanisms and empirical evaluation.

479 Guidelines:

- 480 • The answer [N/A] means that the paper does not include theoretical results.
- 481 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
482 referenced.
- 483 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 484 • The proofs can either appear in the main paper or the supplemental material, but if  
485 they appear in the supplemental material, the authors are encouraged to provide a short  
486 proof sketch to provide intuition.
- 487 • Inversely, any informal proof provided in the core of the paper should be complemented  
488 by formal proofs provided in appendix or supplemental material.
- 489 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 490 4. Experimental result reproducibility

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

494 Answer: [Yes]

495 Justification: The paper describes the evaluated algorithms, datasets, models, serving  
496 backends, hardware platform, and key configuration choices used in the main experiments.  
497 These details provide the information needed to reproduce the experiment.

498 Guidelines:

- 499 • The answer [N/A] means that the paper does not include experiments.
- 500 • If the paper includes experiments, a [No] answer to this question will not be perceived  
501 well by the reviewers: Making the paper reproducible is important, regardless of  
502 whether the code and data are provided or not.
- 503 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
504 to make their results reproducible or verifiable.
- 505 • Depending on the contribution, reproducibility can be accomplished in various ways.  
506 For example, if the contribution is a novel architecture, describing the architecture fully  
507 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
508 be necessary to either make it possible for others to replicate the model with the same  
509 dataset, or provide access to the model. In general, releasing code and data is often  
510 one good way to accomplish this, but reproducibility can also be provided via detailed  
511 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
512 of a large language model), releasing of a model checkpoint, or other means that are  
513 appropriate to the research performed.
- 514 • While NeurIPS does not require releasing code, the conference does require all submis-  
515 sions to provide some reasonable avenue for reproducibility, which may depend on the  
516 nature of the contribution. For example
  - 517 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
518 to reproduce that algorithm.
  - 519 (b) If the contribution is primarily a new model architecture, the paper should describe  
520 the architecture clearly and fully.
  - 521 (c) If the contribution is a new model (e.g., a large language model), then there should  
522 either be a way to access this model for reproducing the results or a way to reproduce  
523 the model (e.g., with an open-source dataset or instructions for how to construct  
524 the dataset).
  - 525 (d) We recognize that reproducibility may be tricky in some cases, in which case  
526 authors are welcome to describe the particular way they provide for reproducibility.  
527 In the case of closed-source models, it may be that access to the model is limited in

528 some way (e.g., to registered users), but it should be possible for other researchers  
529 to have some path to reproducing or verifying the results.

## 530 5. Open access to data and code

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

534 Answer: [Yes]

535 Justification: We provide an anonymized code repo with instructions for setup. Code is  
536 available at <https://anonymous.4open.science/r/FlashEvolve-FEC7/>.

537 Guidelines:

- 538 • The answer [N/A] means that paper does not include experiments requiring code.
- 539 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/  
540 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 541 • While we encourage the release of code and data, we understand that this might not  
542 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not  
543 including code, unless this is central to the contribution (e.g., for a new open-source  
544 benchmark).
- 545 • The instructions should contain the exact command and environment needed to run to  
546 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 547 • The authors should provide instructions on data access and preparation, including how  
548 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 549 • The authors should provide scripts to reproduce all experimental results for the new  
550 proposed method and baselines. If only a subset of experiments are reproducible, they  
551 should state which ones are omitted from the script and why.
- 552 • At submission time, to preserve anonymity, the authors should release anonymized  
553 versions (if applicable).
- 554 • Providing as much information as possible in supplemental material (appended to the  
555 paper) is recommended, but including URLs to data and code is permitted.
- 556

## 557 6. Experimental setting/details

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

560 Answer: [Yes]

561 Justification: The paper specifies the evaluated evolution algorithms, datasets, models,  
562 serving backends, minibatch settings, worker settings, hardware platform, and evaluation  
563 metrics used in the main experiments. Additional details are provided in the code repository.

564 Guidelines:

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

## 570 7. Experiment statistical significance

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

573 Answer: [Yes]

574 Justification: The reported system-efficiency results are obtained from repeated measure-  
575 ments under the same experimental settings, and the tables report mean measured values for  
576 each evaluated workload.

577 Guidelines:

- 578 • The answer [N/A] means that the paper does not include experiments.

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

## 598 8. Experiments compute resources

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

602 Answer: [Yes]

603 Justification: The paper reports the main compute resources used in the experiments,  
604 including the local setup and API-based experiments. It also reports wall-clock time or fixed  
605 time budgets for the main workloads.

606 Guidelines:

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

## 615 9. Code of ethics

616 Question: Does the research conducted in the paper conform, in every respect, with the  
617 NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

618 Answer: [Yes]

619 Justification: The research follows the NeurIPS Code of Ethics. It uses public benchmarks  
620 and existing LLM backends for systems evaluation, and we preserve anonymity in the  
621 submission.

622 Guidelines:

- 623 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of  
624 Ethics.
- 625 • If the authors answer [No], they should explain the special circumstances that require a  
626 deviation from the Code of Ethics.
- 627 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
628 eration due to laws or regulations in their jurisdiction).

## 629 10. Broader impacts

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

632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684

Answer: [Yes]

Justification: The paper discusses that FlashEvolve can reduce the cost of agent evolution and provides a system that is easier to use and deploy.

Guidelines:

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

**11. Safeguards**

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

Answer: [N/A]

Justification: The paper does not release new pretrained models, image generators, or scraped datasets with high misuse risk. It evaluates systems mechanisms using existing models and public benchmark datasets.

Guidelines:

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

**12. Licenses for existing assets**

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

Answer: [Yes]

Justification: The paper cites the existing algorithms, models, serving systems, and benchmark datasets used in the evaluation. We follow their intended research use and provide license and access information in the anonymized artifact.

685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734

Guidelines:

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

**13. New assets**

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

Answer: [Yes]

Justification: The paper releases FlashEvolve as a new code artifact. The anonymized artifact includes documentation, configuration files, and instructions for reproducing the main experiments.

Guidelines:

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

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

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

Answer: [N/A]

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

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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

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

739 Answer: [N/A]

740 Justification: The paper does not involve crowdsourcing experiments or research with human  
741 subjects, so IRB approval or equivalent review is not applicable.

742 Guidelines:

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

#### 753 16. Declaration of LLM usage

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

758 Answer: [Yes]

759 Justification: LLMs are central to the evaluated agent-evolution workflows. The paper  
760 describes the LLM backends used in the experiments, including Qwen3-8B with vLLM and  
761 GPT-4o-mini through API serving.

762 Guidelines:

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