

# CoWork-X: Experience-Optimized Co-Evolution for Multi-Agent Collaboration System

Zexin Lin<sup>\*1,2</sup> Jiachen Yu<sup>\*3</sup> Haoyang Zhang<sup>2</sup> Yuzhao Li<sup>2</sup> Zhonghang Li<sup>1</sup>  
Yujiu Yang<sup>3</sup> Junjie Wang<sup>✉2,3</sup> Xiaoqiang Ji<sup>✉1,4,5</sup>

## Abstract

Large language models are enabling language-conditioned agents in interactive environments, but highly cooperative tasks often impose two simultaneous constraints: sub-second real-time coordination and sustained multi-episode adaptation under a strict online token budget. Existing approaches either rely on frequent in-episode reasoning that induces latency and timing jitter, or deliver post-episode improvements through unstructured text that is difficult to compile into reliable low-cost execution. We propose **CoWork-X**, an active co-evolution framework that casts peer collaboration as a closed-loop optimization problem across episodes, inspired by fast-slow memory separation. CoWork-X instantiates a **Skill-Agent** that executes via HTN (hierarchical task network)-based skill retrieval from a structured, interpretable, and compositional skill library, and a post-episode **Co-Optimizer** that performs patch-style skill consolidation with explicit budget constraints and drift regularization. Experiments in challenging Overcooked-AI-like realtime collaboration benchmarks demonstrate that CoWork-X achieves stable, cumulative performance gains while steadily reducing online latency and token usage.

## 1. Introduction

Large language models (LLMs) accelerate the development of language-conditioned agents in interactive environments and shift communication, planning, and collaboration from

<sup>1</sup>School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen <sup>2</sup>AutoGame Research <sup>3</sup>Tsinghua University <sup>4</sup>School of Artificial Intelligence, The Chinese University of Hong Kong, Shenzhen, <sup>5</sup>Shenzhen Institute of Artificial Intelligence and Robotics for Society. Correspondence to: Junjie Wang <wangjunjie@sz.tsinghua.edu.cn>, Xiaoqiang Ji <jixiaoqiang@cuhk.edu.cn>.

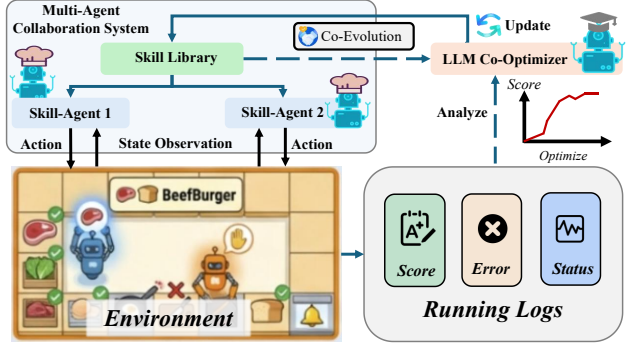


Figure 1. **CoWork-X overview.** Skill-Agents execute via a shared skill library, and an LLM Co-Optimizer updates it from episode logs for closed-loop co-evolution.

handcrafted rules toward a unified learnable paradigm, e.g., in cooperative games. In highly interactive cooperative tasks such as Overcooked-AI, this trend brings two simultaneous requirements. **Real-time coordination** requires stable role assignment and reliable handoffs within a sub-second action loop. **Iterative adaptation** requires continuous strategy updates across repeated trials to handle changes in team partners and environment layouts.

To meet these requirements, prior work broadly progresses along three directions. (i) **In-episode reasoning-driven control:** ReAct (Yao et al., 2023) interleaves reasoning and action to improve interactive decision making. However, in high-frequency interaction, it triggers frequent reasoning calls. These calls introduce online latency and timing jitter, which undermine real-time stability. (ii) **Post-episode language reflection and memory:** Reflexion (Shinn et al., 2023) improves subsequent behavior via after-the-fact reflection on within-episode planning errors. Yet its gains are mainly delivered through textual feedback or prompt updates. Such textual feedback lacks directly callable, structured skill representations and consistency constraints. As a result, experience is not reliably converted into low-cost units for in-episode execution. (iii) **Fast-slow role division, collaboration adaptation, and stronger evaluation pressure:** HLA (Liu et al., 2023b) and DPT-Agent (Zhang et al., 2025) adopt hierarchical or dual-process designs to ease the conflict between real-time execution and reasoning cost. ProAgent (Zhang et al., 2024) emphasizes adaptation

to unfamiliar teammates. Meanwhile, benchmarks such as Collab-Overcooked (Carroll et al., 2019) and OGC (Sarkar et al., 2024) amplify process-level collaboration demands and generalization pressure. This setting exposes a hard tension between continuous adaptation and strict constraints on online budget and stability. In addition, many language-agent paradigms remain *assistant-oriented*, centering on a single LLM serving a human or fixed teammate; fewer studies jointly model all players/co-actors to design a cooperative system that all actors can optimize and co-evolve.

To address these gaps, we propose **CoWork-X**, an active co-evolution framework for peer multi-agent collaboration. CoWork-X shifts the objective from maximizing single-episode performance to achieving *sustained collaboration improvement* across episodes. Motivated by the medial temporal lobe declarative memory system (Squire & Zola-Morgan, 1991), it keeps in-episode behavior fast via rule-like skill memory, and conducts explicit post-episode review to revise and consolidate skills. Accordingly, CoWork-X follows an *Execute–Optimize* closed loop with an implementation-driven structure. (i) **Skill-Agent (HTN-executable skill memory)**: CoWork-X adopts hierarchical task networks (HTN) (Nau et al., 2003) for skill representation, and compiles interpretable, editable, and compositional behaviors into a structured skill library  $\mathcal{S}$ . Within an episode, decision making primarily involves skill retrieval and invocation, reducing reliance on heavyweight online reasoning. (ii) **Co-Optimizer (skill consolidation)**: after each episode, it reads the last trajectory and updates  $\mathcal{S}$  based on outcomes and error signals, revising and consolidating skills for future execution. (iii) **Closed-loop iterative optimization**: across episodes, CoWork-X repeats an execute→diagnose→update→re-execute cycle across episodes. Expensive reasoning and structural edits are pushed into the update stage, keeping in-episode control lightweight while enabling stable gains over time.

We evaluate CoWork-X on Overcooked-AI under settings that stress process-level collaboration and generalization across teammates and layouts. We report cooperative return and online cost (latency, tokens). Across 30 episodes, CoWork-X improves steadily (52.0 at 10 episodes; 96.3 at 30), while ReAct/Reflexion/DPT-WToM achieve 12.5/ – 58.0/3.5 at 10 episodes, respectively. Crucially, CoWork-X executes with 0 online tokens and 2.6 s per episode, approximately  $27\times$  faster than DPT-WToM (71.0 s). These results indicate stable gains from closed-loop iteration with sharply reduced in-episode reasoning burden. Our contributions are:

- We propose **CoWork-X**, an active co-evolution framework for peer multi-agent collaboration that formulates real-time cooperation as a multi-episode closed loop under a strict online budget.
- We introduce an HTN-based **Skill-Agent** with a structured

**skill library** for lightweight in-episode execution, and a **Co-Optimizer** that performs budgeted, regularized patch-style updates for controllable skill evolution.

- We evaluate CoWork-X under strong collaboration and generalization pressure, showing sustained gains from closed-loop iteration with reduced online cost.

## 2. Related Work

### 2.1. Real-Time Interactive Collaboration

Real-time multi-agent coordination under strict latency budgets has inspired diverse architectures. ReAct (Yao et al., 2023) shows that explicit reasoning traces can improve performance, but frequent online LLM calls introduce latency and behavioral jitter that break sub-second coordination. To reduce in-episode cost, hierarchical designs split fast execution from slow deliberation: HLA (Liu et al., 2023b) uses a fast-mind for reactive actions and a slow-mind for high-level reasoning. DPT-Agent (Zhang et al., 2025) further grounds fast control in FSMs with asynchronous reflection and theory-of-mind reasoning. ProAgent (Zhang et al., 2024) infers teammate intent to coordinate with novel partners, but still relies on online LLM reasoning for both action generation and teammate modeling, incurring latency. Despite these advances, prior methods rarely compile experience into structured, verifiable, and compositional skills that support controlled cross-episode evolution, and most remain *assistant-oriented* rather than jointly modeling all actors as peers. Benchmarks such as Overcooked-AI (Carroll et al., 2019) and extensions like OGC (Sarkar et al., 2024) intensify these demands via implicit communication, dynamic roles, and zero-shot transfer, while DPT-Agent (Zhang et al., 2025) further stresses time efficiency. In contrast, we introduce CoWork-X, a peer co-evolution framework that casts coordination as a multi-episode closed loop. It enables stable real-time coordination while consolidating experience into skills that improve controllably across episodes.

### 2.2. Slow Adaptation in Multi-Agent Systems

Sustained multi-agent collaboration requires leveraging past experience to refine coordination over time. Multi-agent RL (e.g., MAPPO (Yu et al., 2022)) can learn effective joint strategies via centralized training with decentralized execution, but typically yields fixed policies that do not support continual self-improvement or adaptation to non-stationary partners. Population-based variants such as COLE (Zhao et al., 2023) broaden partner coverage, yet they still lack post-deployment skill evolution driven by explicit reasoning. Reflection-based agents instead use language feedback for iteration: Reflexion (Shinn et al., 2023) stores textual reflections, and Voyager (Wang et al., 2023) accumulates LLM-generated code skills. However, these updates are largely prompt or memory driven and often miss structural

## CoWork-X: Execute-Optimize Loop

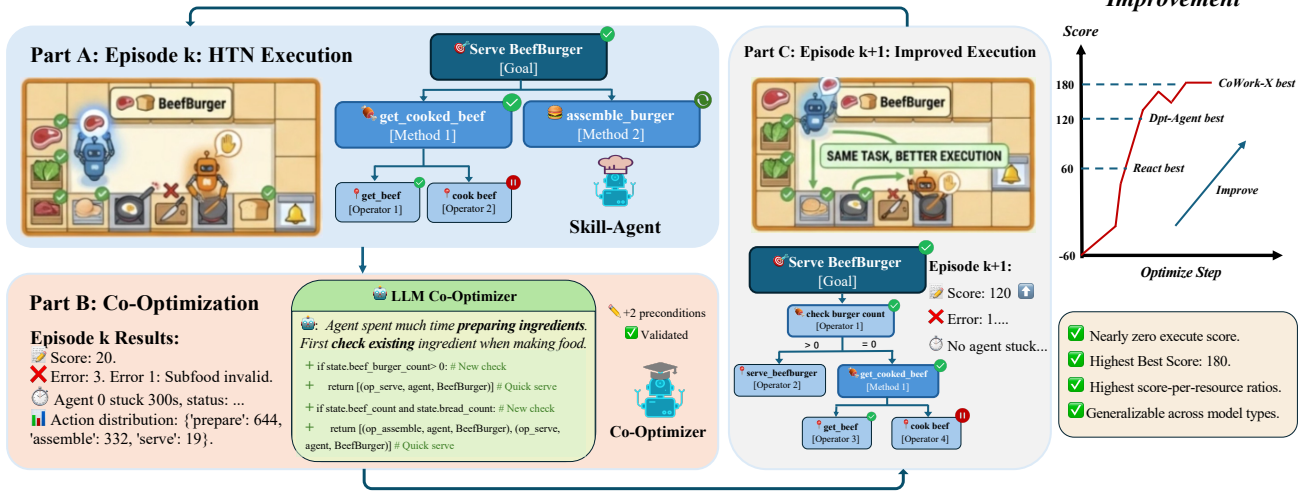


Figure 2. **CoWork-X Execute-Optimize loop.** A Skill-Agent executes an HTN policy from  $S_k$ , then an LLM Co-Optimizer diagnoses episode logs and patches  $S_k \rightarrow S_{k+1}$  (e.g., adding preconditions), improving performance in subsequent episodes.

constraints, verification, and drift regularization needed for controllable accumulation. Classical symbolic planning provides verifiability: HTN (Erol et al., 1996; Nau et al., 2003) decomposes goals into primitive actions, and LLM+P (Liu et al., 2023a) couples LLM understanding with planners for correctness, but the symbolic module is usually treated as a fixed component rather than an evolving repository that closes the loop between execution, diagnosis, and refinement. Meanwhile, most LLM-agent paradigms remain *assistant-oriented*, where a single agent adapts to a human or fixed partner. In contrast, CoWork-X treats all actors as peers sharing a skill library  $S_k$  (iteration- $k$ ), enabling symmetric co-evolution via budgeted post-episode updates under strict online constraints.

### 3. The CoWork-X Framework

#### 3.1. Problem Formulation

We model peer multi-agent collaboration as a fully observable Decentralized MDP (Dec-MDP) with a meta-learning objective over repeated episodes. Formally, the system is a tuple  $\langle \mathcal{I}, \mathcal{X}, \{\mathcal{A}^i\}_{i \in \mathcal{I}}, \mathcal{T}, \mathcal{R}, \gamma \rangle$ , where  $\mathcal{I} = \{1, \dots, N\}$  indexes agents,  $\mathcal{X}$  is the global state, and agent  $i$  selects  $a^i \in \mathcal{A}^i$ . The joint action  $\mathbf{a} \in \mathcal{A} = \prod_{i=1}^N \mathcal{A}^i$  induces transitions  $\mathcal{T} : \mathcal{X} \times \mathcal{A} \rightarrow \Delta(\mathcal{X})$ , and the team reward is given by  $\mathcal{R} : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$  with discount  $\gamma \in [0, 1)$ .

We ground this framework in real-time Overcooked burger preparation scenarios (detailed in Section 4.1) where  $N = 2$  agents collaborate to fulfill time-sensitive food orders. The state space  $\mathcal{X}$  encodes kitchen layout including ingredient stations, cooking equipment, agent positions, and a queue of pending orders with decreasing time-to-deadline for each ac-

tive order. Each agent observes the complete state  $\mathcal{X}$ , though agents must infer teammate intentions from observed actions rather than direct communication. Each agent’s action space  $\mathcal{A}^i$  includes several atomic operations like get beef or cook beef which are defined in low-level game engine architecture. The reward function  $\mathcal{R}$  is controlled by the number of correctly completed orders.

#### 3.2. Overview: Execute-Optimize Closed Loop

As shown in Figure 2, our CoWork-X framework reconciles real-time coordination with cross-episode adaptation via an *Execute-Optimize* loop coupled by a skill library  $S_k$ .  $S_k$  is implemented as a Python file that specifies agent behavior, including state-query utilities, preconditions/procedures for atomic actions, and HTN-encoded task decompositions for diverse orders. In *execution mode*, the Skill-Agent retrieves and invokes skills from  $S_k$  for fast rule-based control. In *optimization mode*, the Co-Optimizer analyzes previous trajectory  $\tau_k$  offline to diagnose recurrent failures and bottlenecks, then updates  $S_k$  by synthesizing patches (e.g., adding missing preconditions or completing recipe logic).

The loop follows  $\text{Execute}(S_k) \rightarrow \text{Diagnose}(\mathcal{D}_k) \rightarrow \text{Update}(S_k \rightarrow S_{k+1}) \rightarrow \text{Execute}(S_{k+1})$ , repeated over iterations. By shifting expensive reasoning to post-hoc offline optimization, CoWork-X keeps online execution lightweight while continuously refining  $S_k$ .

#### 3.3. Skill Agent: HTN-Executable Skill Memory

The agent achieves in-episode execution through an HTN planning framework. The skill library  $S_k$  consists below three hierarchical functions:

- **State representation.** Extracts ingredient quantities from the environment and forms a compact state abstraction that omits spatial details (e.g., object/agent positions). This design prioritizes task-relevant features for fast high-level planning, while delegating navigation and manipulation to the midagent controller.
- **Operators.** Define atomic food-preparation primitives. Preparation operators update ingredient-count attributes to simulate processing; assembly operators first check ingredient availability, then update raw-material and product counters; serving operators validate finished-dish inventory, decrement counts to fulfill orders and accrue rewards, or return error signals when insufficient.
- **Methods.** Map high-level orders to ordered sequences of operators. During execution, a mid-level planner grounds abstract operators into low-level game-engine primitives (e.g., pathfinding and interactions).

### 3.4. Co-Optimizer: Skill Consolidation and Iterative Updates

The skill library is the primary determinant of the Skill-Agent’s game behavior, governing decomposition from high-level objectives to low-level tasks, and is iteratively refined by the Co-Optimizer.

At the onset of the experiment, we initialize  $\mathcal{S}_0$  as a minimal but executable HTN scaffold (Pyhop): a small set of core operator/function templates (e.g., `is_available`, `op_prepare_food`, `op_assemble`, `op_serve`) with conservative placeholder logic, plus a few seed methods (e.g., a basic burger routine) and required registrations. This “infant-state” initialization provides an inductive bias without encoding expert heuristics, while still forcing the Co-Optimizer to infer missing preconditions and state updates, complete recipe logic, and expand task decompositions from execution logs via iterative patches.

After each episode, CoWork-X provides the Co-Optimizer with detailed execution logs beyond the scalar score, including: (i) **Runtime failures** caused by unmet preconditions (e.g., preparing without required ingredients) or invalid variable access. For each failure, we record the exact error location and include the complete environment state and action plans from the two preceding and two succeeding timesteps, which are injected into the prompt to guide repairs to the skill library. (ii) **Stagnation phenomena** flagged when an agent remains inactive for 100 consecutive timesteps. We then provide the environment states and action plans from the first five timesteps of the stagnation interval to help the Co-Optimizer resolve deadlocks. (iii) **Action distribution analysis** reporting the breakdown of action types (e.g., preparation, cooking, serving), which helps identify categories with unusually long durations or low efficiency that warrant optimization. Detailed prompt

used in Co-Optimizer can be found in Appendix B.

We additionally provide the Co-Optimizer with the current skill library file to enforce basic syntactic constraints and discourage radical edits that violate game logic. The best-performing historical library is also included to enable reversion when an update is detrimental. The Co-Optimizer then decides whether further optimization is needed; if so, it leverages these diagnostics to refine rules and outputs an updated Python file as the skill library for the next episode.

## 4. Experimental Settings

### 4.1. Environment and Task Configuration

We evaluate CoWork-X on real-time overcooked-AI-like burger preparation, building on the DPT-Agent environment (Zhang et al., 2025). Unlike the original human-AI setting (one LLM assistant for a human player), we extend it to **symmetric agent-agent collaboration**: both agents are controlled by identical Skill-Agents parameterized by the same skill library  $\mathcal{S}_k$ , instantiating core challenges of Dec-MDP peer coordination.

We adopt Overcooked since it imposes three constraints that directly stress real-time collaboration: (i) *implicit communication*—agents must infer intent from observed actions rather than direct messages; (ii) *tight temporal coupling*—sub-second action loops require low-latency coordination; (iii) *emergent role assignment*—effective play demands dynamic task allocation and handoffs. Together, these properties test whether a method can improve peer coordination consistently across repeated trials.

**Task specification.** Agents prepare and serve three dishes: `BeefBurger` (cooked beef + bread), `LettuceBurger` (chopped lettuce + bread), and `BeefLettuceBurger` (cooked beef + chopped lettuce + bread). Each episode lasts  $T=500$  timesteps with up to four concurrent orders whose time-to-deadline decreases over time for each active order, enforcing prioritization under time pressure. Rewards are +20 per delivered order, +5 per intermediate subtask completion (e.g., cooking beef, chopping lettuce), and −10 for order failure due to timeout or incorrect delivery. Both agents share  $\mathcal{S}_k$ , ensuring performance reflects coordination rather than fixed roles.

**Initial state.** CoWork-X starts from a deliberately impaired but executable skill library: HTN rules are syntactically valid yet semantically incorrect. Operators return unchanged states without precondition checks, and methods decompose tasks without verifying availability. This yields near-zero initial scores and forces the Co-Optimizer to refine coordination through iterative updates. Details of the skill library’s initial state can be found in Appendix D (Figure 9).

**Evaluation protocol.** For each method, we run 30 episodes

under a fixed environment configuration with stochastic order generation. CoWork-X performs 30 Execute–Optimize iterations: episode  $k$  runs with  $\mathcal{S}_k$  to produce  $\tau_k$  and diagnostics, followed by one Co-Optimizer update to obtain  $\mathcal{S}_{k+1}$ . The baselines run 30 episodes independently, without cross-episode skill consolidation. We log per-episode return  $R_k$ , online latency (seconds), and total token usage (prompt + completion across all LLM calls), and report the mean over 30 episodes without outlier removal.

#### 4.2. Baseline Methods

We compare CoWork-X against three state-of-the-art LLM-based multi-agent cooperative-control paradigms.

**ReAct** (Yao et al., 2023) interleaves reasoning and action by generating natural-language traces before issuing executable decisions. We invoke the LLM every 25 timesteps (about 20 calls per  $T=500$  episode) with a rolling event buffer ( $H=1$ ), and parse JSON task assignments for execution. Its main drawback is high online latency from frequent calls and no cross-episode consolidation beyond the prompt window.

**Reflexion** (Shinn et al., 2023) augments ReAct with a reflective layer. A reactive controller runs every 25 timesteps ( $H=1$ ), while a reflection step runs every 75 timesteps ( $H=15$ ) to summarize failures and update textual guidelines appended to subsequent prompts. This adds strategic context but still relies on in-episode LLM calls and does not crystallize reusable skills across episodes.

**DPT-Agent** (Zhang et al., 2025) adopts a dual-process design with asynchronous fast/slow loops and an optional theory-of-mind (ToM) module. The urgent loop triggers every 25 timesteps ( $H=5$ ) for immediate assignments, and the reflection loop every 75 timesteps ( $H=15$ ) for strategic guidance; an FSM fallback covers common subtasks. We evaluate **DPT-WToM** (with ToM) and **DPT-WoToM** (without ToM). Its limitation is reliance on pre-defined structures and prompt-based updates, with no persistent skill evolution beyond the designed scope.

All baselines utilize the same environment, rewards, and LLM backend (Gemini-3-Pro-Preview-Thinking (DeepMind, 2025) via an OpenAI-compatible API), with matched sampling (temperature 0.7, top-p 0.95, max tokens 4096). Differences stem only from control logic and prompting schedules. Detailed Settings are in Appendix C.

#### 4.3. Evaluation Metrics

We evaluate methods along 3 complementary dimensions.

**Task performance.** Episode return  $R_k$  is the total reward over the 500-timestep episode  $k$ , combining order deliveries (+20), intermediate progress (+5), and failures (−10). We report learning curves  $\{R_1, \dots, R_I\}$  over  $I=10$  iterations,

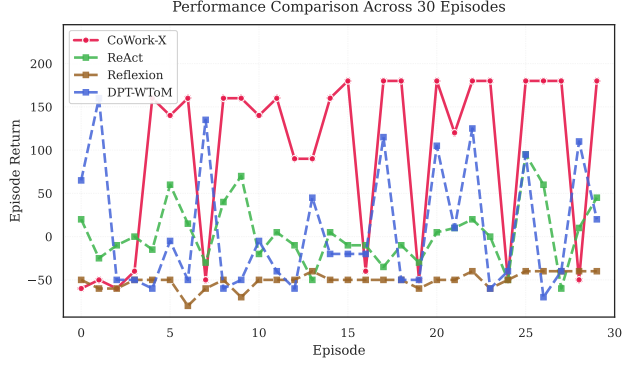


Figure 3. **Performance across 30 episodes.** CoWork-X shows consistent improvement across iterations. Baselines show unstable performance from frequent online LLM calls.

averaged over ten seeds with standard-error bands.

**Online efficiency.** Per-episode latency measures wall-clock runtime including decision overhead (HTN planning for CoWork-X; LLM calls for baselines). Token usage sums prompt+completion tokens across all in-episode LLM calls; CoWork-X uses 0 execution tokens.

**Amortized cost.** Offline optimization tokens count post-episode Co-Optimizer usage. We also report the break-even episode where CoWork-X’s cumulative cost (online + offline) matches baseline cumulative cost (online only).

#### 4.4. Implementation Details

We adopt Gemini-3-Pro-Preview-Thinking as the backbone of the Co-Optimizer. Each trial runs  $I=10$  Execute–Optimize iterations. For each optimization call, we allow up to 3 retries with validation feedback and maintain a history window of the most recent 5 iterations. Stagnation is flagged by 100 consecutive inactive timesteps.

### 5. Results and Analysis

#### 5.1. Main Results

**CoWork-X achieves sustained improvement and remains superior to all baselines throughout the 30-episode trajectory.** Figure 3 shows a clear separation: CoWork-X quickly moves from low initial returns to consistently high performance, and it maintains the top curve for essentially the entire run. Importantly, the curve reflects both (i) a positive learning trend across iterations and (ii) persistent dominance over competing approaches, indicating that improvements are not sporadic but repeatedly recovered and maintained. In contrast, ReAct and DPT-WToM display high-variance swings with frequent regressions, consistent with online LLM-driven control disrupting sub-second coordination (e.g., delayed action selection causing missed

Table 1. Mean performance by episode range (transposed). CoWork-X improves across 10-episode blocks, while baselines remain stagnant or unstable.

Method	0–9	10–19	20–29	Overall (0–29)
ReAct	12.5	-16.5	13.5	3.2
Reflexion	-58.0	-50.0	-45.0	-51.0
DPT-WToM	3.5	-10.5	25.5	6.2
CoWork-X	<b>52.0</b>	<b>109.0</b>	<b>128.0</b>	<b>96.3</b>

Table 2. Computational costs per episode. CoWork-X uses symbolic planning (zero online tokens); optimization costs are amortized across episodes. All measurements report mean values across all episodes.

Method	Online		Total (+ optimization)	
	Time (s)	Tokens	Time (s)	Tokens
ReAct	182.3	79,126	182.3	79,126
Reflexion	67.7	35,635	67.7	35,635
DPT-WToM	71.0	30,090	71.0	30,090
CoWork-X	<b>2.6</b>	<b>0</b>	160.3	22,117

handoffs). Reflexion remains largely negative, suggesting that post-hoc textual reflection alone does not capture the fine-grained spatiotemporal contingencies needed for coordinated cooking/serving under time pressure.

**Block-level aggregation quantifies monotonic gains for CoWork-X and exposes baseline stagnation and instability.** Table 1 confirms progressive refinement: CoWork-X improves from 52.0 (episodes 0–9) to 109.0 (10–19) to 128.0 (20–29), with an overall mean of 96.3. This monotonic rise indicates that cross-episode updates consistently consolidate useful coordination logic into the HTN library rather than overfitting to a single episode. Baselines fail to exhibit a similar pattern. ReAct oscillates around zero (3.2 overall) and even degrades in the middle block (−16.5), aligning with the hypothesis that frequent reasoning calls incur latency/jitter that intermittently breaks coordination. DPT-WToM shows modest late gains (25.5 in 20–29) but remains low on average (6.2 overall), implying limited adaptation beyond its fixed policy scaffold despite added ToM overhead. Reflexion stays strongly negative across all blocks (−58.0, −50.0, −45.0; −51.0 overall), indicating that prompt-level reflections do not reliably translate into executable coordination improvements. Detailed per-episode scores are provided in Appendix A.

## 5.2. Efficiency and Cost Profile

We measure per-episode cost in a 30-episode Overcooked evaluation, reporting (i) *online* runtime and tokens incurred during episode execution, and (ii) *total* cost that additionally includes post-episode optimization.

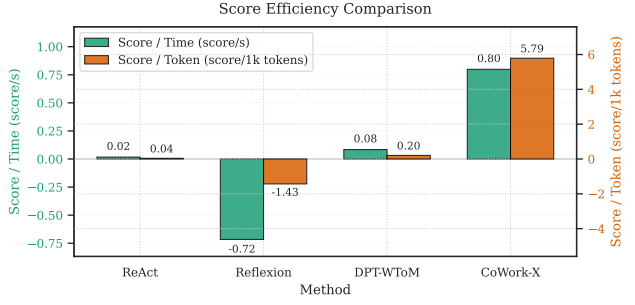


Figure 4. **Score efficiency.** CoWork-X achieves higher score-per-resource ratios: 0.92 score/s and 5.9 score/1k tokens, versus DPT-WToM’s 0.09 and 0.20.

**CoWork-X nearly eliminates online cost, shifting computation to offline optimization.** Table 2 shows that CoWork-X executes with 2.6s online time and 0 online tokens, since the Skill-Agent uses HTN-based symbolic planning without in-episode LLM calls. In contrast, all baselines incur substantial online overhead from frequent LLM invocations, ranging from 67.7–182.3s and 30k–79k tokens per episode.

**Even with optimization included, CoWork-X remains token-efficient while delivering real-time deployment feasibility.** Including post-episode updates, CoWork-X’s total cost is 160.3s and 22,117 tokens per episode on average. Notably, its total token usage is still lower than DPT-WToM (22.1k vs 30.1k), despite paying an explicit optimization budget. Meanwhile, CoWork-X’s online latency is over an order of magnitude smaller than baselines (e.g., 2.6s vs 71.0s for DPT-WToM), making it compatible with sub-second action loops where online LLM calls can systematically break coordination.

## 5.3. Score Efficiency under Online Budgets

We measure *score efficiency* by normalizing average episode return with (i) total wall-clock time and (ii) total token usage, using the same 30-episode evaluation protocol and mean costs/returns aggregated across episodes. Figure 4 shows that CoWork-X attains the highest score/time and score/token ratios, achieving 0.92 score/s and 5.9 score/1k tokens. In comparison, the strongest baseline (DPT-WToM) reaches only 0.09 score/s and 0.20 score/1k tokens, while ReAct is near-zero and Reflexion is negative due to consistently low returns.

**Eliminating online LLM calls is the key driver of real-time efficiency.** CoWork-X’s advantage stems from shifting computation offline: the HTN-based Skill-Agent executes with zero online tokens and minimal latency, while the Co-Optimizer’s cost is amortized across episodes. This separation yields high returns without incurring the online jitter of frequent LLM calls, making CoWork-X compatible with sub-second coordination loops where latency directly trans-

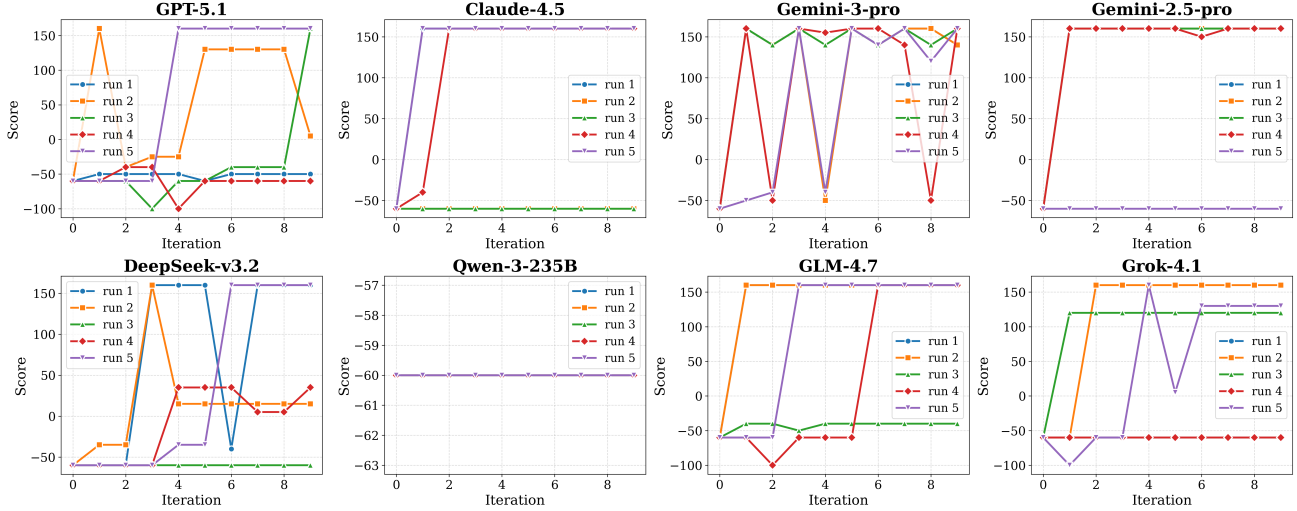


Figure 5. Ablation study of different families of models. Each model was tested in 5 independent runs, each with 10 iterations.

Table 3. Cost comparison at different episode counts. All values report cumulative mean across the first N episodes. CoWork-X’s performance and efficiency advantages increase as optimization costs are amortized.

Range	Method	Time (s)	Tokens	Mean Score	Score Per 1k Tokens
0-9	ReAct	1,823	791,259	12.5	0.016
	Reflexion	677	356,348	-58.0	-0.163
	DPT-WToM	710	300,898	3.5	0.012
	CoWork-X	1,603	221,166	<b>52.0</b>	<b>0.235</b>
0-19	ReAct	3,646	1,582,518	-2.0	-0.001
	Reflexion	1,354	712,695	-54.0	-0.076
	DPT-WToM	1,420	601,795	-3.5	-0.006
	CoWork-X	3,206	442,332	<b>80.5</b>	<b>0.182</b>
0-29	ReAct	5,469	2,373,777	3.2	0.001
	Reflexion	2,031	1,069,043	-51.0	-0.048
	DPT-WToM	2,130	902,693	6.2	0.007
	CoWork-X	4,808	663,499	<b>96.3</b>	<b>0.145</b>

lates into missed handoffs and reduced team reward.

#### 5.4. Cost Amortization and Deployment Scalability

In Table 3, we report cumulative mean statistics over the first  $N \in 10, 20, 30$  episodes (ranges 0–9, 0–19, 0–29). For each method, we aggregate wall-clock time, total tokens, mean score, and score-per-1k-tokens. CoWork-X follows an Execute–Optimize schedule with one post-episode update per episode, whereas baselines incur only in-episode costs.

**CoWork-X maintains a large and early performance lead throughout deployment.** The result shows that CoWork-X is already strong after 10 episodes (mean score 52.0), and the advantage widens with more episodes (80.5 at 20; 96.3 at 30). Baselines remain near zero or negative on average

over the same horizons (e.g., DPT-WToM: 3.5, −3.5, 6.2; Reflexion: −58.0, −54.0, −51.0), indicating little systematic learning from additional trials.

**CoWork-X achieves higher performance at lower cumulative token cost than the strongest baseline.** By 30 episodes, CoWork-X reaches 96.3 using 663k tokens, while DPT-WToM reaches 6.2 using 903k tokens. This yields a  $15.6\times$  score gain with 27% fewer tokens, showing that optimization overhead does not offset overall efficiency.

**Cost amortization improves CoWork-X’s score-per-token profile as episodes accumulate.** CoWork-X’s score-per-1k-tokens decreases from 0.235 (0–9) to 0.182 (0–19) to 0.145 (0–29), consistent with amortizing a fixed optimization cost over more deployments while sustaining high returns. Baselines show no comparable efficiency trend (e.g., DPT-WToM: 0.012  $\rightarrow$  −0.006  $\rightarrow$  0.007), suggesting they primarily accumulate cost without converting experience into reusable coordination capability.

#### 5.5. Generalization across LLM Backbones

We instantiate CoWork-X with diverse LLM backbones spanning 7 model families (GPT (OpenAI, 2025), Gemini (DeepMind, 2025), Claude (Anthropic, 2025), DeepSeek (Liu et al., 2025), Qwen (Yang et al., 2025), GLM (Team et al., 2025), Grok (xAI, 2025)). For each backbone, we run 5 independent trials, each with 10 Execute–Optimize iterations, and report score trajectories in Figure 5.

**CoWork-X generalizes across most backbones, yielding substantial gains from closed-loop optimization.** Across closed-source families, CoWork-X frequently reaches high scores (often  $\geq 140$ ) in multiple runs, indicating that diverse backbones can synthesize and refine executable HTN

## Example of Skill Correction

**Original Rule:** ...if getattr(state, plate\_count) = 0:  
subtasks.append((op\_prepare\_food, agent, Plate))...  
**Error Log:** ...assert food in [Lettuce, Beef, Bread],  
food; **AssertionError: Plate...**  
**Updated Skill:** # Plate availability is assumed im-  
plicitly...delete item Plate.  
**Score:** -40 → 100

Figure 6. An example of skill correction.

## Example of Skill Improvement

**Original Rule:** def is\_available(state, ingredient):  
if ingredient = Bread:  
    return getattr(state, bread\_count) > 0...  
**Error Log:** Agent 1: Failed, lack of necessary ingredi-  
ents to assemble BeefLettuceBurger (ingredients may  
be used by your partner)...  
**Updated Skill:** def is\_available(state, ingredient):  
if ingredient = Bread:  
    # Need 2 to avoid failure when partner grab one.  
    return getattr(state, bread\_count) > 1  
**Score:** 160 → 180

Figure 7. An example of skill improvement.

skills under the same optimization protocol. In contrast, Qwen-3-235B remains near a low-score regime (around -60), revealing a backbone-dependent limit in reliable skill repair and refinement. The failure case can be found in Appendix E.

**Backbones differ in optimization speed and attainable ceilings.** Claude-4.5 and Gemini-2.5-pro typically reach a high ceiling (around 160) within roughly 3 iterations in most runs, whereas GPT-5.1 and DeepSeek-v3.2 improve more slowly and often plateau earlier. GLM-4.7 and Grok-4.1 can match top-end ceilings in some runs, but generally require more iterations to reach a stable optimum.

**Stability varies both within-run and across runs, exposing different failure-recovery profiles.** Some backbones (e.g., Claude-4.5, Gemini-2.5-pro, GLM-4.7) remain stable once a strong library is found, while others (e.g., Gemini-3-pro, GPT-5.1, DeepSeek-v3.2) exhibit larger oscillations, consistent with occasional harmful patches followed by rollback-driven recovery. With “effective optimization” defined as stabilizing above 120 points, the Gemini backbones and GLM-4.7 achieve the highest success rate (4/5), while GPT-5.1 and DeepSeek-v3.2 are less consistent (2/5).

### 5.6. Qualitative Skill Evolution Analysis

We inspect two representative optimization traces produced by the Co-Optimizer (Gemini-3-pro) when updating the HTN-based skill library from episode logs, and visualize the corresponding rule edits, error signals, and score changes

in Figures 6 and 7.

### Log-grounded patches can repair hard execution failures by removing invalid symbolic assumptions.

Figure 6 shows a typical early-stage failure where the Co-Optimizer introduces a non-existent variable (plate\_count) and an unsupported ingredient (Plate), triggering an `AssertionError`. Using the error location and runtime context, the next update deletes the invalid item and restores consistency with the environment’s symbol set, converting a failing trajectory (score -40) into a successful run (score 100). This illustrates that CoWork-X is not merely “prompt tuning”: it performs concrete, verifiable code repairs that directly unblock execution.

### Beyond crash fixes, the Co-Optimizer refines coordination logic by diagnosing implicit interaction bugs.

Figure 7 demonstrates a higher-level improvement where execution does not crash but coordination degrades due to resource contention. The original availability check for bread (bread\_count > 0) allows both agents to compete for the last unit, inducing deadlock and downstream order failure. The Co-Optimizer resolves this by tightening the condition to require bread\_count > 1, effectively encoding a coordination-aware precondition that accounts for teammate consumption. The score increase (from 160 to 180) highlights that CoWork-X can synthesize non-trivial multi-agent constraints from feedback logs even without explicit fault localization.

### Rollback enables recovery from harmful edits and stabilizes multi-step evolution.

In later iterations, the Co-Optimizer may propose overly restrictive rules (e.g., “burger ownership” constraints that prevent efficient handoffs), which can reduce score despite being syntactically valid. CoWork-X mitigates such regressions by allowing the Co-Optimizer to revert to a previously best-performing library, preventing error accumulation from collapsing the optimization trajectory. This failure-recovery pattern complements the quantitative stability trends and is essential for reliable cross-episode co-evolution.

## 6. Conclusion

We presented **CoWork-X**, an active co-evolution framework for peer multi-agent collaboration that reconciles sub-second coordination with cross-episode adaptation under strict online budgets. CoWork-X implements an *Execute-Optimize* loop: Skill-Agents execute lightweight, HTN-structured behaviors from a shared skill library, while an LLM Co-Optimizer performs post-episode diagnosis and patch-style updates from execution logs to consolidate reusable coordination skills. Experiments on real-time overcooked-AI-like benchmark demonstrate sustained improvement and strong efficiency. Across 30 episodes,

CoWork-X steadily increases return (from 52.0 at 10 episodes to 96.3 at 30) while baselines remain low or unstable. By eliminating in-episode LLM calls, CoWork-X achieves 0 online tokens and 2.6s per episode, enabling reliable coordination under tight latency constraints. Overall, CoWork-X shows that structured, log-grounded skill evolution can yield scalable and robust peer collaboration beyond prompt-only adaptation.

## Impact Statements

This work introduces CoWork-X, a closed-loop framework for peer multi-agent collaboration that compiles cross-episode experience into an executable HTN skill library via post-episode LLM-driven updates. The intended impact is to improve coordination reliability under strict online latency budgets, which can benefit real-time interactive systems such as assistive robotics, collaborative agents, and human-in-the-loop decision support.

Potential risks arise from automated skill evolution. Post-episode code updates may introduce brittle heuristics, reward hacking, or silent regressions that only surface under distribution shift (e.g., new layouts or non-stationary partners). LLM-generated patches can also contain implementation bugs or security-relevant flaws, and improved coordination efficiency could be misused to amplify harmful collective behaviors in adversarial settings.

We mitigate these concerns by separating online execution from offline optimization, restricting changes to structured post-episode patches rather than in-episode free-form reasoning, validating generated skill libraries before execution, and keeping a best-performing historical library to support rollback after harmful edits. We also report both performance and cost to make deployment trade-offs explicit. Future work should strengthen safeguards with automated testing and invariant checks, sandboxed execution, explicit safety constraints during optimization, and broader evaluation under adversarial and out-of-distribution conditions.

To support reproducibility and responsible reuse, we release our implementation under the MIT License. We encourage users to follow the license terms (e.g., attribution and preserving the license notice). We also encourage publicly sharing derivatives and results built on CoWork-X to improve transparency, comparability, and collective progress.

## References

- Anthropic. Claude sonnet 4.5 system card. System card, Anthropic, October 2025. URL <https://www.anthropic.com/claude-sonnet-4-5-system-card>.
- Carroll, M., Shah, R., Ho, M. K., Griffiths, T., Seshia, S.,

Abbeel, P., and Dragan, A. On the utility of learning about humans for human-ai coordination. In *Advances in Neural Information Processing Systems*, volume 32, 2019.

DeepMind, G. Gemini 3 pro model card. Technical report, Google, 2025. URL <https://storage.googleapis.com/deepmind-media/Model-Cards/Gemini-3-Pro-Model-Card.pdf>.

Erol, K., Hendler, J., and Nau, D. S. Umcp: A sound and complete procedure for hierarchical task-network planning. *Artificial Intelligence*, 94(1-2):249–292, 1996.

Liu, A., Mei, A., Lin, B., Xue, B., Wang, B., Xu, B., Wu, B., Zhang, B., Lin, C., Dong, C., et al. Deepseek-v3. 2: Pushing the frontier of open large language models. *arXiv preprint arXiv:2512.02556*, 2025.

Liu, B., Jiang, Y., Zhang, X., Liu, Q., Zhang, S., Biswas, J., and Stone, P. Llm+p: Empowering large language models with optimal planning proficiency. *arXiv preprint*, arXiv:2304.11477, 2023a. URL <https://arxiv.org/abs/2304.11477>.

Liu, J., Liu, C., Zhou, Y., Zhao, Y., and Chen, L. Llm-powered hierarchical language agent for real-time human-ai coordination. *arXiv preprint*, arXiv:2312.15224, 2023b. URL <https://arxiv.org/abs/2312.15224>.

Nau, D. S., Au, T.-C., Ilghami, O., Kuter, U., Murdock, J. W., Wu, D., and Yaman, F. SHOP2: An HTN planning system. *Journal of Artificial Intelligence Research*, 20: 379–404, 2003.

OpenAI. Gpt-5 system card. Technical report, OpenAI, 2025. URL <https://cdn.openai.com/gpt-5-system-card.pdf>.

Sarkar, A., Rahman, M., Ciosek, K., Murray, S., Griffin, C., Everett, T., Lewis, R., Graham, M., Hendriks, V., Campbell, J., Slumbers, O., et al. The overcooked generalisation challenge. *arXiv preprint*, arXiv:2406.17949, 2024. URL <https://arxiv.org/abs/2406.17949>.

Shinn, N., Cassano, F., Gopinath, A., Narasimhan, K., and Yao, S. Reflexion: Language agents with verbal reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 36, pp. 8634–8652, 2023.

Squire, L. R. and Zola-Morgan, S. The medial temporal lobe memory system. *Science*, 253(5026):1380–1386, 1991.

Team, V., Hong, W., Yu, W., Gu, X., Wang, G., Gan, G., Tang, H., Cheng, J., Qi, J., Ji, J., Pan, L., Duan, S., Wang, W., Wang, Y., Cheng, Y., He, Z., Su, Z., Yang, Z., Pan, Z.,

- Zeng, A., Wang, B., Chen, B., Shi, B., Pang, C., Zhang, C., Yin, D., Yang, F., Chen, G., Xu, J., Zhu, J., Chen, J., Chen, J., Chen, J., Lin, J., Wang, J., Chen, J., Lei, L., Gong, L., Pan, L., Liu, M., Xu, M., Zhang, M., Zheng, Q., Yang, S., Zhong, S., Huang, S., Zhao, S., Xue, S., Tu, S., Meng, S., Zhang, T., Luo, T., Hao, T., Tong, T., Li, W., Jia, W., Liu, X., Zhang, X., Lyu, X., Fan, X., Huang, X., Wang, Y., Xue, Y., Wang, Y., Wang, Y., An, Y., Du, Y., Shi, Y., Huang, Y., Niu, Y., Wang, Y., Yue, Y., Li, Y., Zhang, Y., Wang, Y., Wang, Y., Zhang, Y., Xue, Z., Hou, Z., Du, Z., Wang, Z., Zhang, P., Liu, D., Xu, B., Li, J., Huang, M., Dong, Y., and Tang, J. Glm-4.5v and glm-4.1v-thinking: Towards versatile multimodal reasoning with scalable reinforcement learning, 2025. URL <https://arxiv.org/abs/2507.01006>.
- Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., Fan, L., and Anandkumar, A. Voyager: An open-ended embodied agent with large language models. *Transactions on Machine Learning Research*, 2023.
- xAI. Grok-4 model card. Model card, xAI, August 2025. URL <https://data.x.ai/2025-08-20-grok-4-model-card.pdf>.
- Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., and Cao, Y. React: Synergizing reasoning and acting in language models. In *Proc. of the 11th Int. Conf. on Learning Representations (ICLR 2023)*, 2023. Oral presentation; also arXiv:2210.03629.
- Yu, C., Velu, A., Vinitisky, E., Gao, J., Wang, Y., Bayen, A., and Wu, Y. The surprising effectiveness of PPO in cooperative multi-agent games. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. URL <https://openreview.net/forum?id=YVXaxB6L2Pl>.
- Zhang, C., Yang, K., Hu, S., Wang, Z., Li, G., Sun, Y., Zhu, C., Zhong, Z., Zhang, W., Ouyang, W., et al. Proagent: Building proactive cooperative agents with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 19547–19555, 2024.
- Zhang, S., Wang, X., Zhang, W., Li, C., Song, J., Li, T., Qiu, L., Cao, X., Cai, X., Yao, W., Zhang, W., Wang, X., and Wen, Y. Leveraging dual process theory in language agent framework for real-time simultaneous human-AI collaboration. In Che, W., Nabende, J., Shutova, E., and Pilehvar, M. T. (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4081–4108, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.206. URL <https://aclanthology.org/2025.acl-long.206/>.
- Zhao, Y., Gao, Y., Wu, J., Yuan, L., Wang, G., Huang, K., Jiang, Y., Yang, Y., Wu, Y., and Du, Y. Tackling cooperative incompatibility for zero-shot human-ai coordination. *arXiv preprint, arXiv:2306.03034*, 2023. URL <https://arxiv.org/abs/2306.03034>.

Table 4. Episode-by-episode scores for single experimental run. CoWork-X shows consistent improvement across iterations while baselines remain stagnant.

Episode	ReAct	Reflexion	DPT-WToM	CoWork-X
0	20	-50	65	-60
1	-25	-60	160	-50
2	-10	-60	-50	-60
3	0	-50	-50	-40
4	-15	-50	-60	160
5	60	-50	-5	140
6	15	-80	-50	160
7	-30	-60	135	-50
8	40	-50	-60	160
9	70	-70	-50	160
10	-20	-50	-5	140
11	5	-50	-40	160
12	-10	-50	-60	90
13	-50	-40	45	90
14	5	-50	-20	160
15	-10	-50	-20	180
16	-10	-50	-20	-40
17	-35	-50	115	180
18	-10	-50	-50	180
19	-30	-60	-50	-50
20	5	-50	105	180
21	10	-50	10	120
22	20	-50	125	180
23	0	-40	-60	180
24	-50	-60	-40	-50
25	95	-50	95	180
26	60	-40	-70	180
27	-60	-40	-40	180
28	10	-40	110	-50
29	45	-40	20	180

## Appendix

### A. Raw Experimental Results

Table 4 presents the complete episode-by-episode scores from the single experimental run of 30 episodes for each method. CoWork-X demonstrates progressive improvement through iterative skill library refinement, while baselines show stagnant or unstable performance due to lack of cross-episode learning mechanisms.

#### Key observations:

- **CoWork-X progressive improvement:** Scores increase from 52.0 (episodes 0–9) to 109.0 (episodes 10–19) to 128.0 (episodes 20–29), demonstrating effective skill library optimization across iterations. Summary statistics are presented in the main text (Table 1).
- **Baseline stagnation:** ReAct, Reflexion, and DPT-WToM show no consistent improvement across episode ranges, as they lack cross-episode learning mechanisms.
- **Final convergence:** CoWork-X achieves 180 points in

7 of the final 10 episodes (episodes 20–29), indicating convergence to effective coordination strategies.

- **Occasional failures:** CoWork-X shows occasional negative scores (-50 to -60) even in later episodes, typically occurring when the Co-Optimizer attempts exploratory modifications that temporarily degrade performance before recovery in subsequent iterations.

### B. Prompt of Co-Optimizer

Figure 8 shows the prompt used by Co-Optimizer. The prompt context consists the current skill library, the origin library, historical performance metrics, and diagnostic error logs. It requires the Co-Optimizer to refine the skill library by addressing the specific runtime errors identified in the previous iteration, with the final output strictly constrained to Python format.

### C. Model Specifications

Table 5. LLM model names and their corresponding version identifiers.

Model Name	Version Identifier
<i>Main Experiments:</i>	
Gemini-3-Pro	gemini-3-pro-preview-11-2025-thinking
<i>Model Ablation Study:</i>	
GPT-5.1	gpt-5.1-all
Claude-4.5	claude-sonnet-4-5-20250929
Gemini-3-Pro	gemini-3-pro-preview-11-2025-thinking
Gemini-2.5-Pro	gemini-2.5-pro
DeepSeek-v3.2	deepseek-v3.2
Qwen-3-235B	qwen3-235b-a22b-instruct-2507
GLM-4.7	glm-4.7
Grok-4.1	grok-4.1

**Implementation details.** All models use identical sampling parameters: temperature 0.7, top-p 0.95, and max tokens 4096. The main experiments (ReAct, Reflexion, DPT-WToM, and CoWork-X baselines) use Gemini-3-Pro-Preview-Thinking exclusively to ensure fair comparison. The model ablation study (Section 5.3) evaluates CoWork-X across different model families to demonstrate generalizability.

### D. Details of Skill Library

**Original Skill Library.** Figure 9 illustrates the initial, minimal yet executable skill library. This baseline comprises solely fundamental state representations, function definition for methods and operators, and their corresponding Pyhop declarations. Within the method functions, we provide exemplar subtask lists to prevent the Co-Optimizer from generating logic incompatible with the game environment. For the remaining functions, we supply default state

returns accompanied by instructional comments to guide the Co-Optimizer in function completion. Under this original rule set, Gemini-3-pro achieved a score of -60. We designate this infancy state as the initialization point for the CoWork-X optimization process.

**Best-Optimized Skill Library.** Figure 10 illustrates the optimal skill library derived from the 18th optimization iteration of Gemini-3-pro. Guided by this library, the skill agent attained a peak score of 180. Notably, subsequent iterations achieving this score utilized structurally similar libraries, indicating that the optimization process effectively converged at this point. This library provides comprehensive and precise definitions for all methods and operators, ensuring that the agent executes the game without triggering any runtime errors.

## E. Failure Mode

Figure 11 illustrates an instance of optimization failure exhibited by Qwen-3-225B. Constrained by inherent model capabilities, Qwen-3-225B accumulated numerous erroneous variable definitions—such as `beef_fresh_count` across multiple execution epochs. Within the game environment, the model failed to accurately pinpoint the error locations based on the provided error logs, thereby impeding the Co-Optimizer’s effective self-optimization.

## F. License

The source code will be available under the MIT License. The full text of the license is provided at <https://opensource.org/licenses/MIT>.

## Prompt of Co-Optimizer

You are an expert in HTN (Hierarchical Task Network) planning and Python programming.  
 We are optimizing a rule file for an autonomous agent in an Overcooked-like environment using the Pyhop planner.  
 {history\_context}  
 {best\_rules\_context}

Current Rule File {current\_rule\_file}:  
 {current\_rules}

origin Rule File {origin\_rule\_file}:  
 {origin\_rules}

Notice: Every variant name and function name in origin rule file is correct, although origin rule is broken. Do not change any basic elements in origin rule file, only to complete, optimize the given function and write new process logic of other recipes.

Experiment Result (Iteration {iteration}):  
 Score: {score}  
 Log Analysis: {analysis log}

Task:

Analyze the logs and the current rules to identify if the rule is perfect.

Common issues include:

1. **\*\*Precondition Failures\*\***: Operators returning 'False' because state checks fail.
2. **\*\*State Update Failures\*\***: Operators not updating the state correctly (e.g., not incrementing counts).
3. **\*\*Method Logic\*\***: Methods not decomposing into the correct sequence of operators.
4. **\*\*Missing Logic\*\***: Handling for specific ingredients or burgers is missing.
5. **\*\*Recipe Errors\*\***: Incorrect ingredient counts (e.g., requiring 2 items when 1 is standard).

Your Goal: Modify the rule file to FIX the errors (if exist) and MAXIMIZE the score.

- If you see some recipe logic missing, add it.
- If you see "Assemble" failing, check the ingredient requirements and counts.
- If you see a method missing a step, add it.
- Learn from the optimization history - don't repeat mistakes from previous iterations.

CRITICAL INSTRUCTIONS:

- Output ONLY valid Python code.
- Do NOT include any explanations, markdown formatting, or comments outside the code.
- Do NOT start with text like 'Looking at the code...' or 'Here is the fix...'.
- The first line MUST be valid Python (import statement, comment, or function definition)
- Your entire response will be saved directly as a .py file
- If you think the rule is already perfect, output current rules totally the same.

Start your response with the first line of Python code immediately.

Figure 8. Prompt of Co-Optimizer.

## Original Skill Library

```

from pyhop import hop as pyhop

# determine if some ingredient is available. You need to complete the logic corresponding to this function
according to the variable name given in the log.
def is_available(state, ingredient):
    return True

# prepare food and update state. You need to complete the update logic corresponding to this function according
to the variable name given in the log.
def op_prepare_food(state, agent, food_type):
    return state

# assemble used ingredient and update state. You need to complete the update logic corresponding to this
function according to the variable name given in the log.
def op_assemble(state, agent, burger_type):
    return state

# serve finished burger and update state. You need to complete the update logic corresponding to this function
according to the variable name given in the log.
def op_serve(state, agent, burger_type):
    return state

# return subtasks can only contains functions start with op_, and end with assemble and serve. do not use
recursion
def m_make_beef_burger(state, agent):
    if getattr(state, 'beef_burger_count', 0) > 0:
        return [('op_serve', agent, 'BeefBurger')]
    if is_available(state, 'Beef'):
        return [('op_assemble', agent, 'BeefBurger'), ('op_serve', agent, 'BeefBurger')]
    return False

# declare defined methods and operators
def declare_rules(pyhop_module):
    pyhop_module.declare_operators(op_prepare_food, op_assemble, op_serve)
    pyhop_module.declare_methods('make_BeefBurger', m_make_beef_burger)

```

Figure 9. Original Skill Library.

## Optimized Best Skill Library

```

from pyhop import hop as pyhop
# Determine if some ingredient is available. Buffer 'Bread' (> 1) to mitigate race conditions.
def is_available(state, ingredient):
    if ingredient == 'Beef': return getattr(state, 'beef_cooked_count', 0) > 0
    elif ingredient == 'Lettuce': return getattr(state, 'lettuce_chopped_count', 0) > 0
    elif ingredient == 'Bread': return getattr(state, 'bread_count', 0) > 1 # Require buffer
    return False

# Prepare food and update state counters.
def op_prepare_food(state, agent, food_type):
    if food_type == 'Beef': state.beef_cooked_count = getattr(state, 'beef_cooked_count', 0) + 1
    elif food_type == 'Lettuce': state.lettuce_chopped_count = getattr(state, 'lettuce_chopped_count', 0) + 1
    elif food_type == 'Bread': state.bread_count = getattr(state, 'bread_count', 0) + 1
    return state

# Assemble used ingredients and update state.
def op_assemble(state, agent, burger_type):
    bread = getattr(state, 'bread_count', 0)
    if bread < 1: return False
    if burger_type == 'BeefBurger':
        beef = getattr(state, 'beef_cooked_count', 0)
        if beef >= 1:
            state.bread_count -= 1; state.beef_cooked_count -= 1;
            state.beef_burger_count = getattr(state, 'beef_burger_count', 0) + 1;
            return state
    elif burger_type == 'LettuceBurger':
        lettuce = getattr(state, 'lettuce_chopped_count', 0)
        if lettuce >= 1:
            state.bread_count -= 1; state.lettuce_chopped_count -= 1;
            state.lettuce_burger_count = getattr(state, 'lettuce_burger_count', 0) + 1;
            return state
    elif burger_type == 'BeefLettuceBurger':
        beef, lettuce = getattr(state, 'beef_cooked_count', 0), getattr(state, 'lettuce_chopped_count', 0)
        if beef >= 1 and lettuce >= 1:
            state.bread_count -= 1; state.beef_cooked_count -= 1;
            state.lettuce_chopped_count -= 1;
            state.beef_lettuce_burger_count = getattr(state, 'beef_lettuce_burger_count', 0) + 1;
            return state
    return False

# Serve finished burger and update state.
def op_serve(state, agent, burger_type):
    if burger_type == 'BeefBurger':
        if getattr(state, 'beef_burger_count', 0) > 0: state.beef_burger_count -= 1; return state
    elif burger_type == 'LettuceBurger':
        if getattr(state, 'lettuce_burger_count', 0) > 0: state.lettuce_burger_count -= 1; return state
    elif burger_type == 'BeefLettuceBurger':
        if getattr(state, 'beef_lettuce_burger_count', 0) > 0: state.beef_lettuce_burger_count -= 1; return state
    return False

# Method: Make BeefBurger
def m_make_beef_burger(state, agent):
    if getattr(state, 'beef_burger_count', 0) > 0: return [('op_serve', agent, 'BeefBurger')]
    tasks = []
    if not is_available(state, 'Beef'): tasks.append(('op_prepare_food', agent, 'Beef'))
    if not is_available(state, 'Bread'): tasks.append(('op_prepare_food', agent, 'Bread'))
    return tasks + [('op_assemble', agent, 'BeefBurger'), ('op_serve', agent, 'BeefBurger')]

# Method: Make LettuceBurger
def m_make_lettuce_burger(state, agent):
    if getattr(state, 'lettuce_burger_count', 0) > 0: return [('op_serve', agent, 'LettuceBurger')]
    tasks = []
    if not is_available(state, 'Lettuce'): tasks.append(('op_prepare_food', agent, 'Lettuce'))
    if not is_available(state, 'Bread'): tasks.append(('op_prepare_food', agent, 'Bread'))
    return tasks + [('op_assemble', agent, 'LettuceBurger'), ('op_serve', agent, 'LettuceBurger')]

# Method: Make BeefLettuceBurger (Priority: Beef -> Lettuce -> Bread)
def m_make_beef_lettuce_burger(state, agent):
    if getattr(state, 'beef_lettuce_burger_count', 0) > 0: return [('op_serve', agent, 'BeefLettuceBurger')]
    tasks = []
    if not is_available(state, 'Beef'): tasks.append(('op_prepare_food', agent, 'Beef'))
    if not is_available(state, 'Lettuce'): tasks.append(('op_prepare_food', agent, 'Lettuce'))
    if not is_available(state, 'Bread'): tasks.append(('op_prepare_food', agent, 'Bread'))
    return tasks + [('op_assemble', agent, 'BeefLettuceBurger'), ('op_serve', agent, 'BeefLettuceBurger')]

# Declare defined methods and operators
def declare_rules(pyhop_module):
    pyhop_module.declare_operators(op_prepare_food, op_assemble, op_serve)
    pyhop_module.declare_methods('make_BeefBurger', m_make_beef_burger)
    pyhop_module.declare_methods('make_LettuceBurger', m_make_lettuce_burger)
    pyhop_module.declare_methods('make_BeefLettuceBurger', m_make_beef_lettuce_burger)

```

Figure 10. Best-optimized skill library by Gemini-3-pro.

## Failure Mode

```

from pyhop import hop as pyhop
# Determine if some ingredient is available. Logic based on variable names from logs.
def is_available(state, ingredient):
    if ingredient == 'Beef': return (getattr(state, 'beef_fresh_count', 0) > 0 or getattr(state, 'beef_in_progress_count', 0) > 0 or
    getattr(state, 'beef_cooked_count', 0) > 0)
    elif ingredient == 'Lettuce': return (getattr(state, 'lettuce_unchopped_count', 0) > 0 or getattr(state, 'lettuce_chopped_count',
    0) > 0)
    elif ingredient == 'Bread': return getattr(state, 'bread_count', 0) > 0
    return False
# Prepare food and update state.
def op_prepare_food(state, agent, food_type):
    if food_type == 'Beef':
        if getattr(state, 'beef_fresh_count', 0) > 0: state.beef_fresh_count -= 1; state.beef_in_progress_count = getattr(state,
        'beef_in_progress_count', 0) + 1; return state
        return False
    elif food_type == 'Lettuce':
        if getattr(state, 'lettuce_unchopped_count', 0) > 0: state.lettuce_unchopped_count -= 1; state.lettuce_chopped_count =
        getattr(state, 'lettuce_chopped_count', 0) + 1; return state
        return False
    elif food_type == 'Bread': state.bread_count = getattr(state, 'bread_count', 0) + 1; return state
    return False
# Cook beef (in progress -> cooked) & Chop lettuce (unchopped -> chopped)
def op_cook_beef(state, agent):
    if getattr(state, 'beef_in_progress_count', 0) > 0: state.beef_in_progress_count -= 1; state.beef_cooked_count = getattr(state,
    'beef_cooked_count', 0) + 1; return state
    return False
def op_chop_lettuce(state, agent):
    if getattr(state, 'lettuce_unchopped_count', 0) > 0: state.lettuce_unchopped_count -= 1; state.lettuce_chopped_count =
    getattr(state, 'lettuce_chopped_count', 0) + 1; return state
    return False
# Assemble used ingredient and update state.
def op_assemble(state, agent, burger_type):
    if burger_type == 'BeefBurger':
        if getattr(state, 'beef_cooked_count', 0) < 1 or getattr(state, 'bread_count', 0) < 1: return False
        state.beef_burger_count = getattr(state, 'beef_burger_count', 0) + 1; state.beef_cooked_count -= 1; state.bread_count -= 1;
    return state
    elif burger_type == 'LettuceBurger':
        if getattr(state, 'lettuce_chopped_count', 0) < 1 or getattr(state, 'bread_count', 0) < 1: return False
        state.lettuce_burger_count = getattr(state, 'lettuce_burger_count', 0) + 1; state.lettuce_chopped_count -= 1;
    state.bread_count -= 1; return state
    elif burger_type == 'BeefLettuceBurger':
        if getattr(state, 'beef_cooked_count', 0) < 1 or getattr(state, 'lettuce_chopped_count', 0) < 1 or getattr(state, 'bread_count',
    0) < 1: return False
        state.beef_lettuce_burger_count = getattr(state, 'beef_lettuce_burger_count', 0) + 1; state.beef_cooked_count -= 1;
    state.lettuce_chopped_count -= 1; state.bread_count -= 1; return state
    return False
def op_serve(state, agent, burger_type):
    if burger_type == 'BeefLettuceBurger' and getattr(state, 'beef_lettuce_burger_count', 0) > 0:
        state.beef_lettuce_burger_count -= 1; state.burger_count = getattr(state, 'burger_count', 0) + 1; state.total_score += 15;
    return state
    elif burger_type == 'BeefBurger' and getattr(state, 'beef_burger_count', 0) > 0:
        state.beef_burger_count -= 1; state.burger_count = getattr(state, 'burger_count', 0) + 1; state.total_score += 10; return
    state
    elif burger_type == 'LettuceBurger' and getattr(state, 'lettuce_burger_count', 0) > 0:
        state.lettuce_burger_count -= 1; state.burger_count = getattr(state, 'burger_count', 0) + 1; state.total_score += 10; return
    state
    return False
def m_make_beef_burger(state, agent):
    if getattr(state, 'beef_burger_count', 0) > 0: return [(('op_serve', agent, 'BeefBurger'))]
    if is_available(state, 'Beef') and is_available(state, 'Bread'):
        return [(('op_prepare_food', agent, 'Beef'), ('op_cook_beef', agent), ('op_prepare_food', agent, 'Bread'), ('op_assemble', agent,
        'BeefBurger'), ('op_serve', agent, 'BeefBurger'))]
    return False
def m_make_lettuce_burger(state, agent):
    if getattr(state, 'lettuce_burger_count', 0) > 0: return [(('op_serve', agent, 'LettuceBurger'))]
    if is_available(state, 'Lettuce') and is_available(state, 'Bread'):
        return [(('op_prepare_food', agent, 'Lettuce'), ('op_chop_lettuce', agent), ('op_prepare_food', agent, 'Bread'), ('op_assemble',
        agent, 'LettuceBurger'), ('op_serve', agent, 'LettuceBurger'))]
    return False
def m_make_beef_lettuce_burger(state, agent):
    if getattr(state, 'beef_lettuce_burger_count', 0) > 0: return [(('op_serve', agent, 'BeefLettuceBurger'))]
    if is_available(state, 'Beef') and is_available(state, 'Lettuce') and is_available(state, 'Bread'):
        return [(('op_prepare_food', agent, 'Beef'), ('op_cook_beef', agent), ('op_prepare_food', agent, 'Lettuce'), ('op_chop_lettuce',
        agent), ('op_prepare_food', agent, 'Bread'), ('op_assemble', agent, 'BeefLettuceBurger'), ('op_serve', agent, 'BeefLettuceBurger'))]
    return False
def declare_rules(pyhop_module):
    pyhop_module.declare_operators(op_prepare_food, op_cook_beef, op_chop_lettuce, op_assemble, op_serve)
    pyhop_module.declare_methods('make_BeefBurger', m_make_beef_burger)
    pyhop_module.declare_methods('make_LettuceBurger', m_make_lettuce_burger)
    pyhop_module.declare_methods('make_BeefLettuceBurger', m_make_beef_lettuce_burger)

```

Figure 11. A failed skill library in runs of Qwen-3-225B.