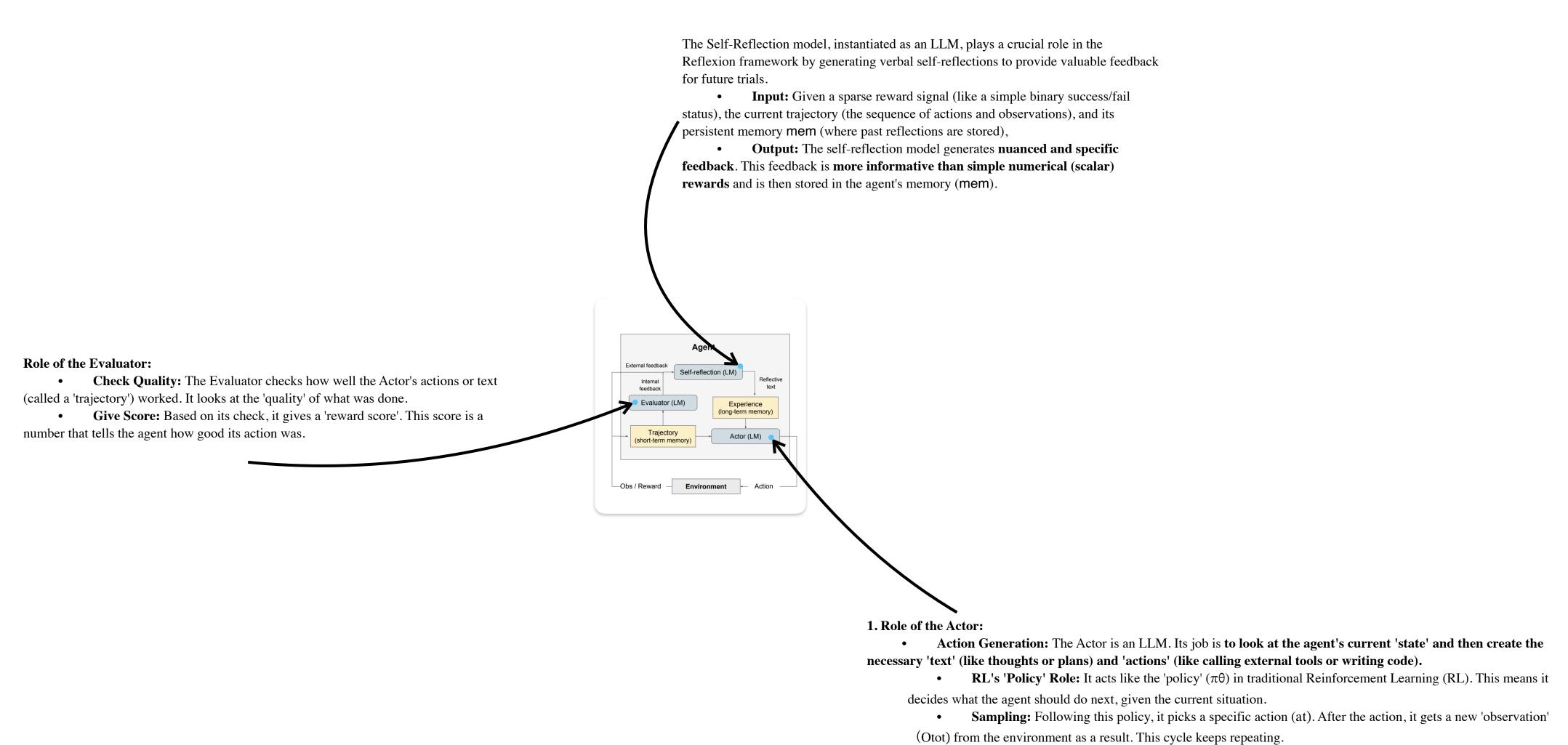


I --> J{Increment Trial t = t+1};

J --> C;



AlfWorld experiment for Reflexion

AlfWorld Environment: • Text-based, multi-step decision-making tasks: (e.g., finding, moving, manipulating objects). • Self-evaluation is essential: The environment only indicates task completion, requiring the agent to self-evaluate its progress. • Action Generator: ReAct from Yao et al. [30] was used to generate actions.

• Natural Language Classification using LLM: The LLM directly evaluates trajectories and provides feedback. • Hand-written Heuristics: • Triggers self-reflection if the same action/response repeats more than 3 times (indicating inefficiency or "hallucination"). • Triggers self-reflection if the number of actions exceeds 30 (indicating inefficient planning).

• Role of the Self-Reflection Module: This module analyzes failed trajectories and reward information to generate specific verbal feedback (Sr^t) about what went wrong and why, not just that it was wrong. This feedback acts as a crucial "self-hint" to prevent the agent from repeating the same mistakes in future attempts.

Comparison: Reflexion vs. Baseline (ReAct-only): • Baseline (ReAct-only): If self-reflection is deemed necessary, the environment is reset, and a new trial begins without the self-reflection process. • Reflexion Applied: If self-reflection is needed, the agent analyzes its mistakes through self-reflection, updates its memory, then resets the environment and starts a new trial.

Memory Management: • Strategic Memory Storage: To stay within LLM context limits, the memory stores only the last 3 self-reflections (experiences). This is a strategic choice to keep the most relevant and useful information. • Contribution to Long-term Learning: This memory allows the agent to effectively use past lessons, maintain long-term context, and anticipate similar errors in complex problem-solving. • Initial Setup: Domain-specific few-shot examples were provided to prevent syntactic errors in the initial agent's behavior.

• Significant Performance Boost: ReAct + Reflexion successfully completed 130 out of 134 AlfWorld tasks, showing markedly superior performance compared to the ReAct-only approach. • Error Detection & Handling: Reflexion effectively detected and managed "hallucinations" and inefficient planning using its heuristics. • Continuous Learning: ReAct + Reflexion demonstrated the ability to learn and solve additional tasks over 12 consecutive trials. • Limitations of Baseline: In contrast, the ReAct-only approach showed a performance plateau, with improvement stopping between trials 6 and 7. The "hallucination" issues and lack of long-term improvement in ReAct-only strongly highlight the need for self-correction mechanisms like Reflexion.

(a) ALFWorld Success Rate ReAct only - hallucination - ReAct only - inefficient planning ReAct + Reflexion (Heuristic) -- ReAct + Reflexion (GPT) ReAct + Reflexion - hallucination ReAct + Reflexion - inefficient planning Figure 3: (a) AlfWorld performance across 134 tasks showing cumulative proportions of solved tasks

using self-evaluation techniques of (Heuristic) and (GPT) for binary classification. (b) Classification

of AlfWorld trajectories by reason of failure.

Reasoning: HotPotQA

Reasoning-Only Setup:

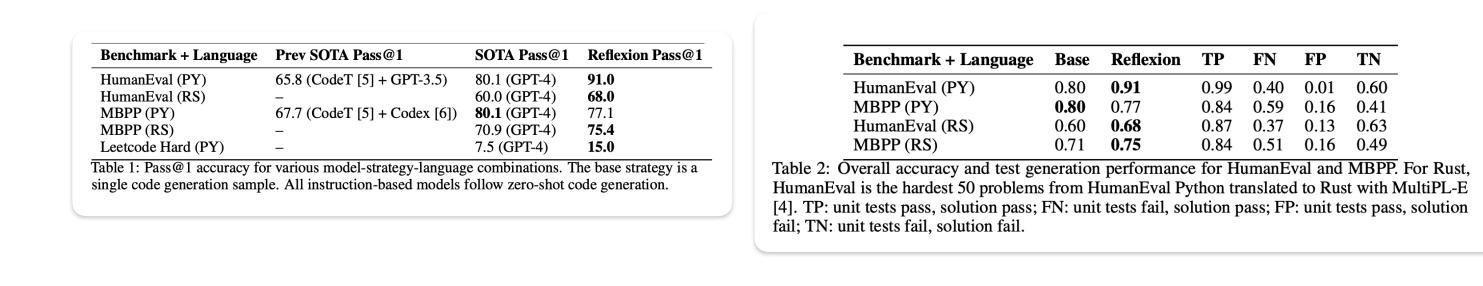
• Used **Reflexion** + **Chain-of-Thought** (**CoT**). • Agent was given the ground truth context (Cgt) to isolate 6-shot prompting for CoT. Holistic Q&A Setup: • Used **Reflexion** + **ReAct**. • Agent could retrieve context using a Wikipedia API and infer answers with step-by-step thinking. • 2-shot prompting for ReAct. • Self-Reflection Implementation: • 2-shot prompting for self-reflection. • Used **exact match answer grading** for a binary success signal between trials. • Self-reflection loop amplified this binary signal, similar to AlfWorld, with a memory size of **3 experiences**. • Reflexion agents retried failed tasks until 3 consecutive failures on that task. Key Results: • Reflexion significantly outperformed all baseline approaches (ReAct-only, CoT-only, CoT (GT)-only) across multiple learning steps. • Baseline approaches **failed to probabilistically improve** on any tasks; once failed, they couldn't solve them later. • Even CoT (GT) (with ground truth context) failed on 39% of questions. • Reflexion improved accuracy by 14% for CoT (GT) by helping the agent correct mistakes without needing access to the ground truth answer during the correction process. • This shows Reflexion's ability to **enhance reasoning and self-correction** even in complex Q&A tasks. **№** 0.8 (a) HotPotQA Success Rate (c) HotPotQA Episodic Memory - ◆- · CoT only → CoT (GT) + Reflexion ReAct only □ 0.9 --- CoT (GT) EPM ● 0.6 ← CoT + Reflexion → CoT (GT) EPM + Reflexion ReAct + Reflexion ້ 0.4 \rightarrow ____ Figure 4: Chain-of-Thought (CoT) and ReAct. Reflexion improves search, information retrieval, and reasoning capabilities on 100 HotPotQA questions. (a) Reflexion ReAct vs Reflexion CoT (b) Reflexion CoT (GT) for reasoning only (c) Reflexion vs episodic memory ablation.

Analysis

• Goal: To isolate the benefit of the self-reflective step for reasoning, using CoT (GT) as the baseline. • CoT (GT) Baseline: Tests reasoning ability over long contexts by using Chain-of-Thought with provided ground truth context. **Ablation Steps:**

Baseline: CoT (GT) only. Episodic Memory (EPM) added: Includes the most recent trajectory. **Reflexion Agent:** Implements the standard self-reflection step as a final pass.

Intuition Tested: Whether agents learn more effectively by iteratively using verbal explanations (written in the first person). • Key Finding (from Figure 4): Self-reflection improves learning by an 8% absolute boost over the episodic memory learning advantage. • Conclusion: This result supports the argument that "refinement-only" approaches are less effective than "self-reflection-guided refinement" approaches.



Conclusion

Reflexion, an approach that leverages verbal reinforcement to teach agents to learn from past mistakes. We empirically show that Reflexion agents significantly outperform currently widely-used decision-making approaches by utilizing self-reflection. In future work, Reflexion could be used to employ more advanced techniques that have been thoroughly studied in traditional RL settings, such as value learning in natural language or off-policy exploration techniques.

Sample "Reflection"code

• **Ma:** Actor (An LLM that generates actions). • Me: Evaluator (An LLM or external system that assesses the outcome or trajectory of the generated actions). • Msr: Self-Reflection Model (An LLM that generates reflection **mem:** The agent's self-reflection memory (in the form of a list).

• max_trials: The maximum number of attempts/trials. • **Me.pass**(): Indicates whether the Evaluator signifies a 'success' state (e.g., passing a coding problem).

import random # Used for simulated random outcomes in the example # --- Define Mock/Abstract Components ---

In a real implementation, this would involve prompting an LLM.

"""Represents the Actor LLM (Ma) that generates actions.""" def __init__(self, Ilm_parameters=None): self.llm_parameters = llm_parameters # Simulates LLM configuration (e.g., model name, temperature) def generate_action(self, state: str, memory: list) -> str: Generates the next action based on the current state and the agent's memory.

print(f"Actor: Current state is '{state}'. Reflecting on memory: {memory}") # Example: Generate a hypothetical action string # The policy $\pi\theta$ (ai |si) is parameterized by Ma and the memory (mem) action = f"Attempt based on current state '{state}' and " \ f"past reflection '{memory[-1] if memory else 'no previous insights'}'" return action class Evaluator: """Represents the Evaluator (Me) that assesses the outcome of a trajectory."""

def __init__(self): self._passed = False # Initial state: task not passed self.trial_count = 0 def evaluate_trajectory(self, trajectory: str) -> str: Evaluates the given trajectory and returns a feedback string. In a real scenario, this would involve interacting with an external environment (e.g., running code, checking game state) and analyzing its output. self.trial_count += 1

print(f"Evaluator: Evaluating trajectory: '{trajectory}'") # Simulated evaluation logic: # 50% chance to pass randomly, or higher chance after a few trials, # or if trajectory explicitly contains "correct" (for demonstration) if "correct" in trajectory.lower() or (self.trial_count > 2 and random.random() < 0.7): self._passed = True feedback = "SUCCESS: The task was completed correctly."

feedback = f"FAILURE: The task failed. Reason: {random.choice(['syntax error', 'logic error', 'timeout', 'incorrect output'])}. Trajectory was: '{trajectory}'" return feedback def is_passed(self) -> bool: """Returns True if the task has been successfully passed, False otherwise."""

return self._passed class SelfReflectionModel: """Represents the Self-Reflection Model (Msr) that generates reflections."""

def __init__(self): def generate_reflection(self, feedback: str, trajectory: str, memory: list) -> str: Generates a self-reflection text based on the feedback, the executed trajectory, and potentially the past memory. In a real implementation, this would involve prompting an LLM with the feedback. print(f"SelfReflectionModel: Generating reflection for feedback: '{feedback}'") if "SUCCESS" in feedback: reflection = f"Successfully completed the task. Key learning from this success: {feedback}. Previous trajectory: '{trajectory}'" # For failures, suggest actionable insights reflection = f"Failed. Need to re-evaluate approach. Feedback: {feedback}. Previous trajectory: '{trajectory}'. " \ f"Consider trying a different strategy or debugging {random.choice(['variable names', 'loop conditions', 'API calls', 'overall plan'])}."

return reflection # --- Reflexion Algorithm Implementation -def run_reflexion_agent(max_trials: int = 5): Executes the learning loop for a Reflexion agent as described in the pseudocode. # Initialize Actor, Evaluator, Self-Reflection: Ma, Me, Msr ma = Actor(Ilm_parameters={"model": "GPT-4", "temperature": 0.7}) # Ma me = Evaluator() msr = SelfReflectionModel()

Initialize policy $\pi\theta$ (ai |si), θ = {Ma, mem} # The policy is implicitly defined by how the Actor (Ma) uses the current state and memory. mem = [] # Agent's episodic memory buffer, initialized as empty # Generate initial trajectory using $\pi\theta$ print("\n--- Initial Trial (t=0) ---") current_state = "Starting the task from scratch."

The first action is generated without any prior reflection in memory initial_action = ma.generate_action(current_state, mem) # A simple representation of the trajectory (sequence of actions/observations) initial_trajectory = f"Initial attempt: '{initial_action}'" # Evaluate τ0 using Me initial_feedback = me.evaluate_trajectory(initial_trajectory) # Generate initial self-reflection sr0 using Msr initial_reflection = msr.generate_reflection(initial_feedback, initial_trajectory, mem) # Set mem ← [sr0] mem.append(initial_reflection)

print(f"Memory after initial reflection: {mem}\n")

Set t = 0 (loop counter)

while Me not pass or t < max trials do # Continue trials as long as the task is not passed AND we haven't exceeded max_trials while not me.is_passed() and t < max_trials: t += 1 # Increment trial counter for the current iteration

Generate τt = [a0, o0, . . . ai , oi] using $\pi \theta$ # The Actor now uses the accumulated 'mem' (reflections) as part of its policy current_state = f"Continuing the task based on previous attempts. Trial count: {t}" action = ma.generate_action(current_state, mem) # In a real system, 'action' would be executed in the environment, # and 'observation' would be received, forming the full trajectory. # Here, we simplify the trajectory representation. trajectory = f"Action in trial {t}: '{action}'" # Check if task was already passed in a previous iteration (edge case) if me.is_passed(): print("Task already passed, no new action needed for this iteration.")

break # Exit the loop if already passed # Evaluate τt using Me feedback = me.evaluate_trajectory(trajectory) # Generate self-reflection srt using Msr reflection = msr.generate_reflection(feedback, trajectory, mem) # Append srt to mem mem.append(reflection) print(f"Memory after trial {t} reflection: {mem}\n") # end while

return (Final results summary) print("\n--- Final Results ---") if me.is_passed(): print(f"Task completed successfully after {t} trials!") print(f"Task not completed within {max_trials} trials.") print(f"Final Memory: {mem}") # Display the last feedback received, or initial if only one trial occurred print(f"Final Feedback: {initial_feedback if t==0 else feedback}") # Execute the Reflexion agent simulation if __name__ == "__main__": run_reflexion_agent(max_trials=5)