

Try-Then-Eval: Equipping an LLM-based Agent with a Two-Phase Mechanism to Solve Computer Tasks

Vy Le^{1,2}

¹ University of Information Technology, Ho Chi Minh City, Vietnam

² Vietnam National University, Ho Chi Minh City, Vietnam

What ?

We develop an **autonomous LLM-based agent** capable of performing **web automation tasks** directly from **natural language task descriptions**, without relying on human demonstrations.

Main contributions:

- Propose a **two-phase Try-Then-Eval mechanism** for computer task automation.
- Eliminate the need for **manual expert demonstrations**.
- Enable the agent to **learn from its own successes and failures**.

Why ?

- Web automation tasks are **common but repetitive**, ranging from clicking buttons to filling complex forms.
- Existing approaches often **depend heavily on handcrafted demonstrations**.
- Large Language Models (LLMs) offer strong reasoning ability but **lack structured self-improvement mechanisms**.

Our approach enables **scalable, adaptive, and demonstration-free automation**.

Overview

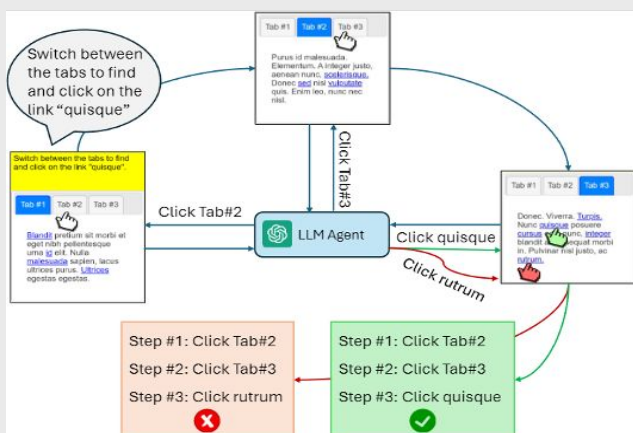


Figure 1. Example of computer task performed by the LLMs agent.

The agent operates in a closed interaction loop:

1. Observe the current web interface (shortened DOM).
2. Generate candidate actions using an LLM.
3. Execute actions and observe outcomes.
4. Store successful trials and extract rules from failures.
5. Reuse learned experience to improve future decisions.

This process allows the agent to self-improve over time.

Description

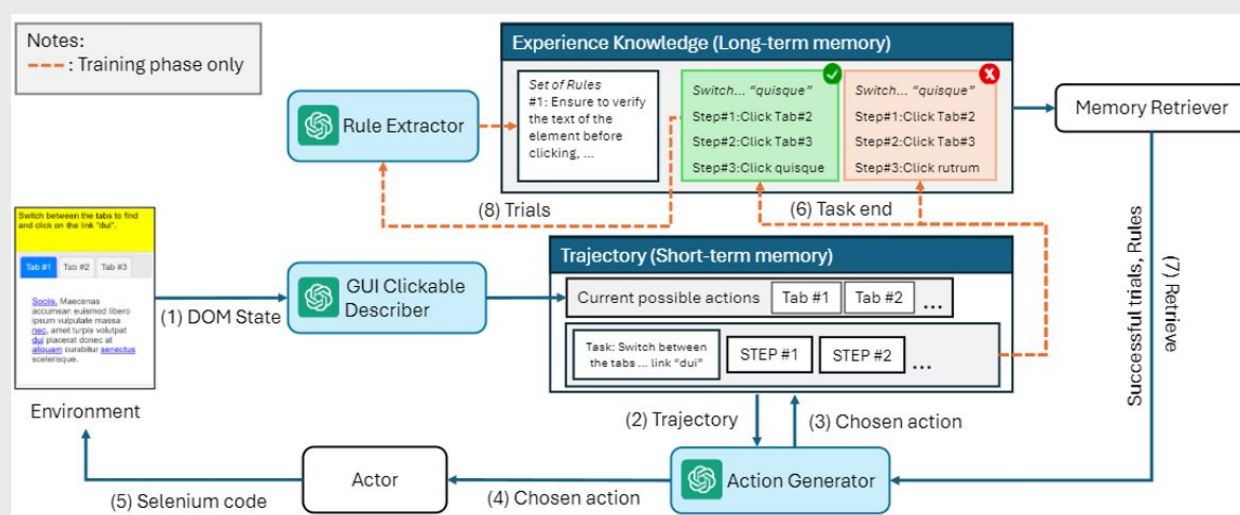


Figure 2. Overview of proposed method.

1. Shortened DOM Representation

- Only visible and interactive elements are extracted.
- Reduces noise and token usage.
- Converts DOM elements into natural language descriptions.

2. Iterative Planning

At each step, the agent:

- Considers task goal.
- Reviews past actions (short-term memory).
- Selects the most suitable action.

Planning continues until task success or failure.

3. Experience Reinforcement

Training phase:

- Successful trials are stored.
- Failed trials are used to extract textual rules.

Evaluation phase:

- Agent reuses the best rules and trials.
- No new rules are added.

This mechanism enables learning without human supervision.

4. Dataset and metric

Dataset: Experiments are conducted on the **MiniWoB++** benchmark, consisting of diverse web-based interaction tasks.

Metric: Performance is evaluated using **Success Rate (SR)**, defined as the ratio of successfully completed tasks.