# M1 Internship Defense

DynPol : Dynamic Policy Library for web agents

Martin CUINGNET
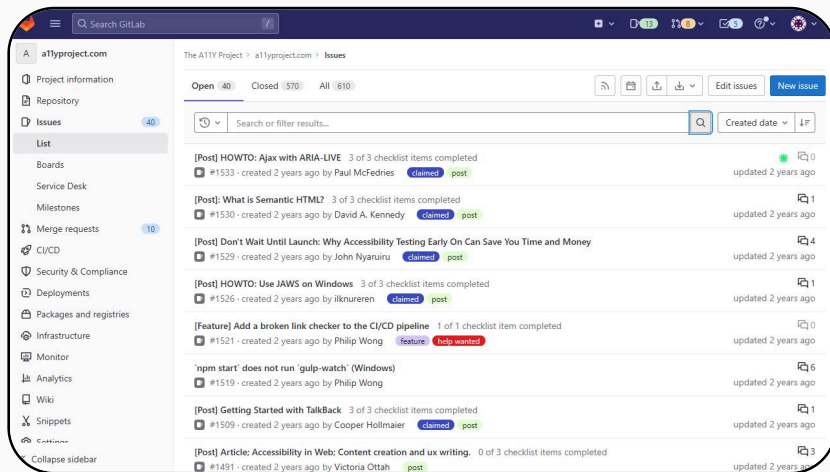
Aldo LIPANI, Jerome RAMOS, Bin WU | Web Intelligence Team, UCL

2025-12-05

ENS Paris-Saclay

# 1. Context

# 1.1 What is a web agent ?



Figure 1: Web Environment

$\Rightarrow$

Web

Agent

$\Rightarrow$

**Action**

- click [380]

- type [67] [Carnegie Mellon]

Example of web task :

**Find the highest rated book on the Shopping site**

**Sort issues based on a label in a GitLab repo**

The agent needs to:
- Find the associated repo
- Go to the issues page
- Sort issues based on the label
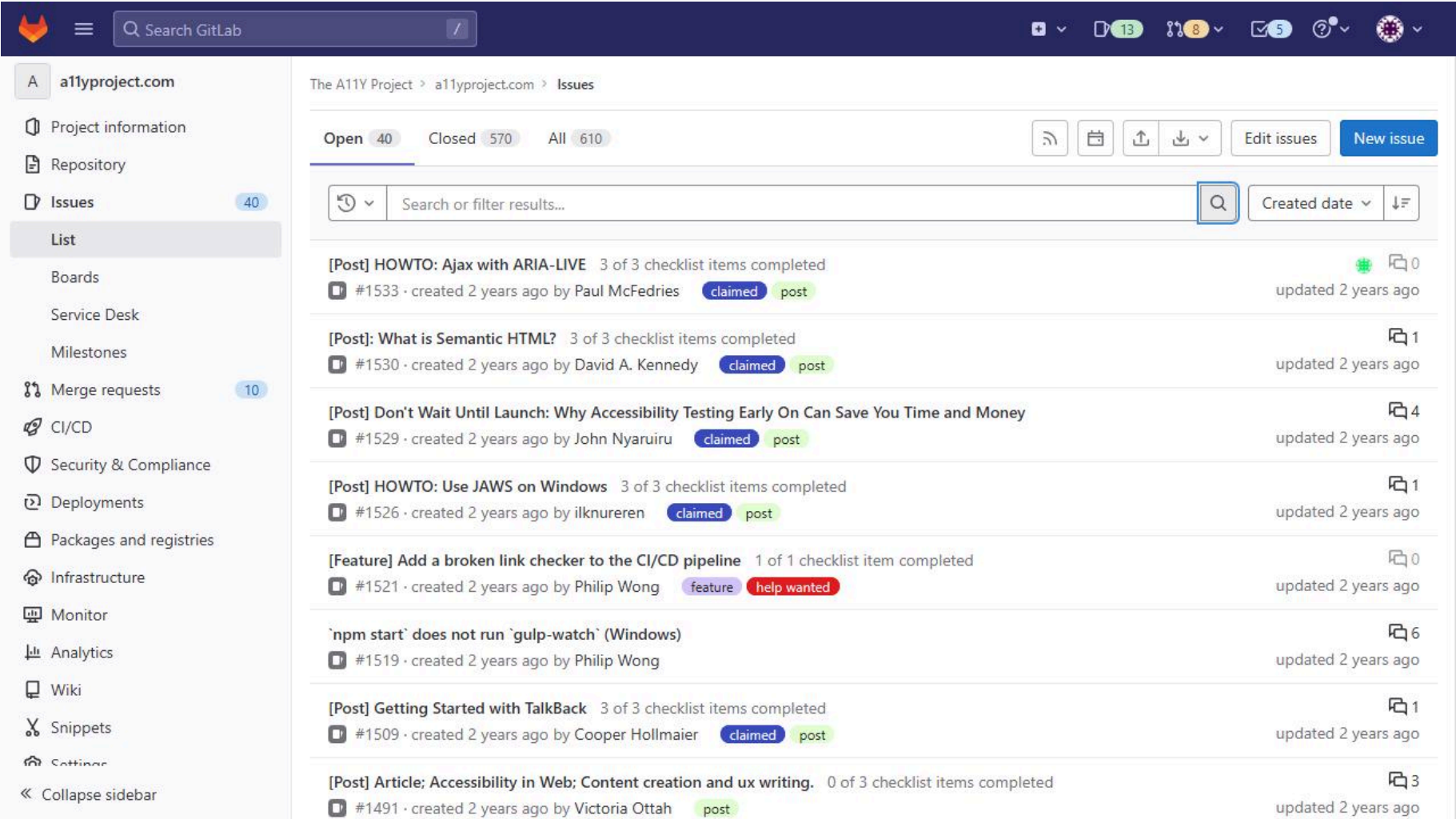
# 1.2 Web Navigation Example



Figure 2: Issue page before sorting
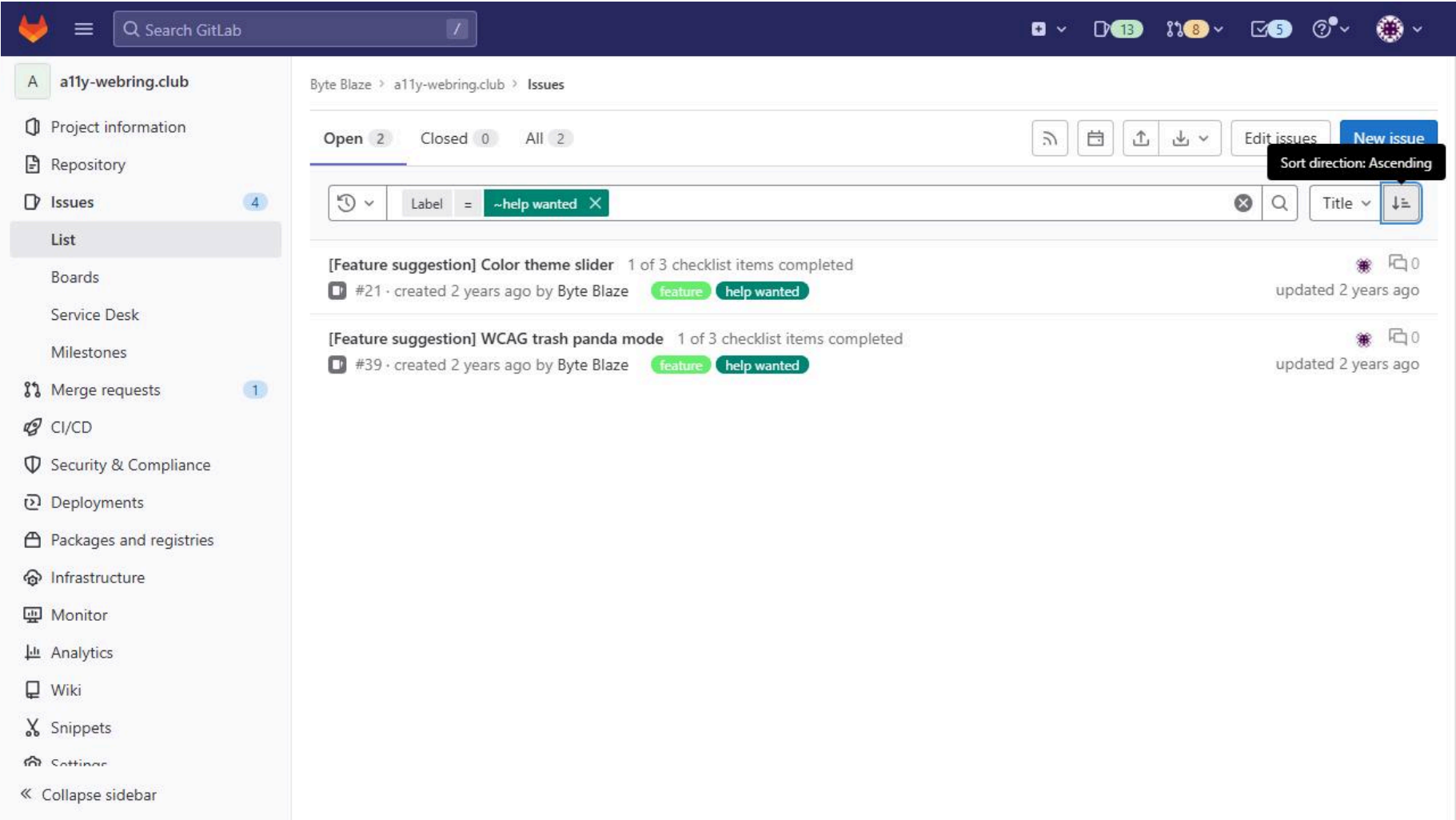
# 1.2 Web Navigation Example



Figure 3: Issue page after successfully sorting by clicking on the `help wanted` tag

# 1.3 Task Examples

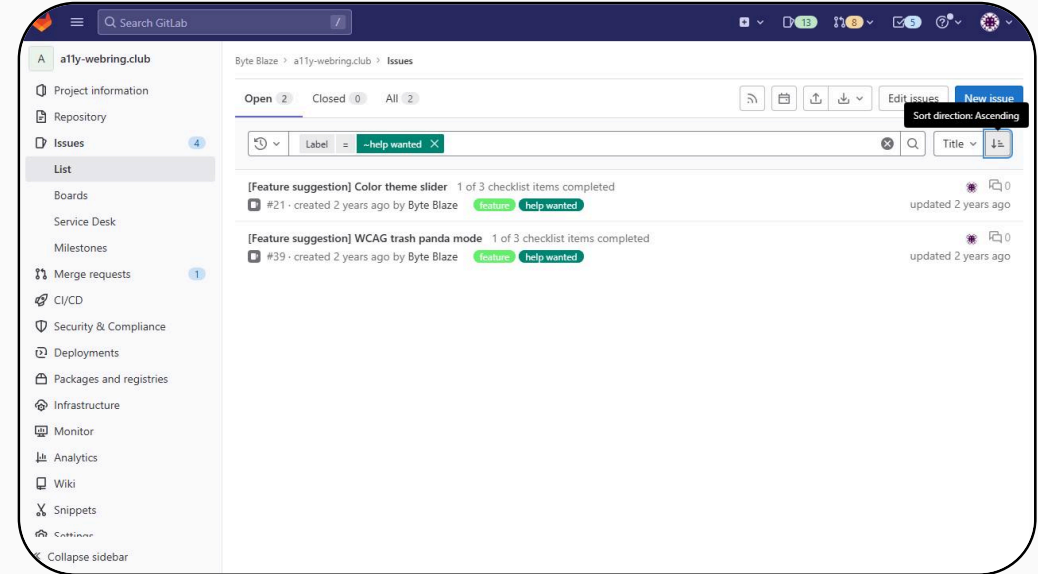| Task ID | Website | Task Objective |
|---------|---------|----------------|
| 181 | Gitlab | Open my latest created issue that has theme editor in its title to check if it is closed |
| 298 | Shopping Admin | Show the most recent completed order |
| 386 | Shopping | What is the rating of Ugreen lightning to 3.5mm cable. Please round to the nearest whole number |
| 643 | Reddit | Post a notice on a virtual meetup for racing cars enthusiasts on Oct 21st in the nyc subreddit |
| 767 | Maps | Find the walkway to the closest chain grocessory owned by a local business from 401 Shady Ave, Pittsburgh. |

# 1.4 Challenges with tasks in web environments

## Ambiguity with web structures

- Complex objects

- Variety of websites

- Action effect unclear

## Complex tasks composed of multiple subtasks

- Complex tasks can be divided in multiple, simpler subtasks

- Hard to keep track of what to do without being lost



- Find the associated repo
- Go to the issues page
- Sort issues based on the label

**SteP** (Stacked LLM Policies for Web Actions) [1] : a web-agent paper

**Proposed solutions** :

- Complex tasks : Simpler subtasks

- Ambiguity : External knowledge

Combine the two solutions $\Rightarrow$ **policies**

Policy : a **strategy** to solve a **specific subtask**.

- Contains **instructions** to solve the subtask
- **Called** by the agent and handled **independently**
- Works like a **non-deterministic subroutine**

Example of policy :

`filter_and_sort_issues`

Click the button to reveal the sorting options.\n*   Select the desired sorting criteria from the dropdown menu.\n*   Check the current sort direction. If the current sort direction is not the desired one, click the sort direction button to toggle it.\nPlease, use ONLY page operation and no subroutine actions.",

In the previous example the policy can be called with :
`filter_and_sort_issues [help wanted]`

SteP policies: **handcrafted in advance**

**Problems** :

- **Static** Policy Library : No adaptability
- **Handcrafted** policies : Misalignment

Solution : Make the policy library **dynamic**

Two components to do so :

- **Self-Improvement Mechanism** : Experience $\Rightarrow$ Policy instruction optimization

- **Automatic Curriculum** : Planning + Create new policies

# 2. What is **DynPol** ?

# 2.1 Policy Stack

**Policy Stack**

```
_____

  find_issue
_____

  find_repo
```
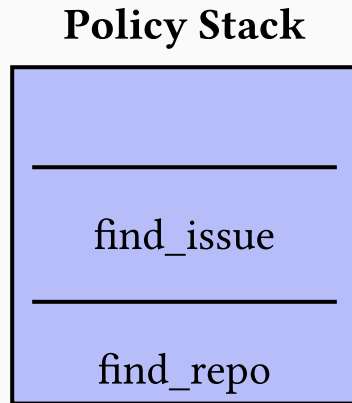
Figure 4: Policy Stack

Policy **currently active** = Policy on **top of the stack**

Calling a policy $\Rightarrow$ adding policy on top of the stack

Finishing a policy $\Rightarrow$ popping top policy of the stack

**Dynamic Policy Library**

create_group
find_subreddit
find_issue
search_customer
find_repo

Figure 5: Policy Library

Store all **available** policies.

Two ways to edit it:
- **Self-Improvement Mechanism** $\Rightarrow$ Improve policy instruction
- **Automatic Curriculum** $\Rightarrow$ Create new policy

Relevant
Policies

Figure 6: Relevant Policies

$|\text{LLM Context Size}| < \infty \Rightarrow$ Fixed number of policies passed to agent

How to decide which policies to fetch ?
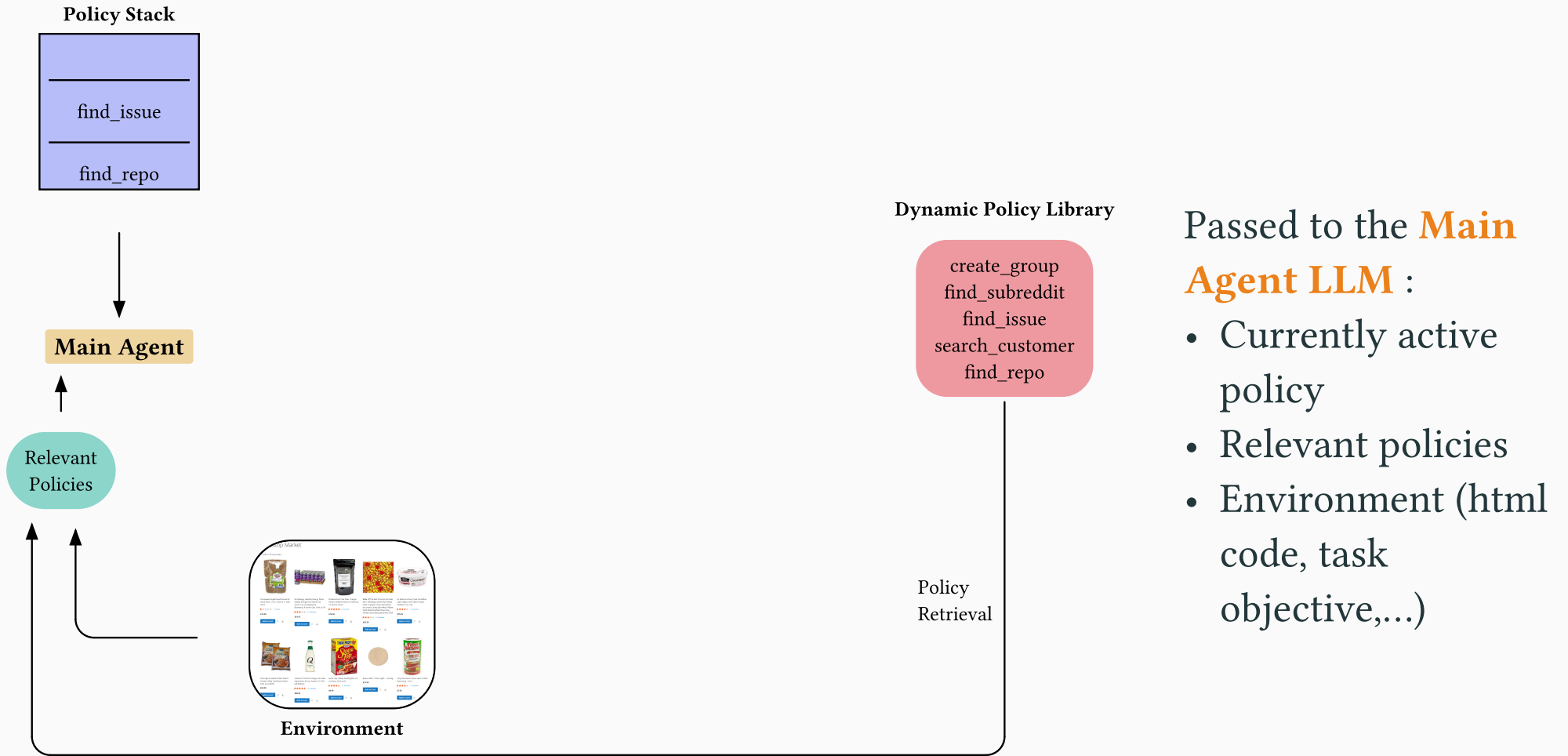
$\mapsto$ Policy embedding closest to Task embedding

**Policy Stack**

| |
|---|
| |
| find_issue |
| find_repo |

**Main Agent**

Relevant Policies

**Environment**

**Dynamic Policy Library**

create_group
find_subreddit
find_issue
search_customer
find_repo

Policy Retrieval

Passed to the **Main Agent LLM** :
- Currently active policy
- Relevant policies
- Environment (html code, task objective,...)

Figure 7: DynPol Main Loop - Step 1

**Policy Stack**

**find_issue**

**find_repo**

**Dynamic Policy Library**

create_group
find_subreddit
find_issue
search_customer
find_repo

**Main Agent**

outputs → Action

Relevant Policies

**Environment**

Policy Retrieval

**Main Agent** issues actions to solve the **task objective**

Figure 8: DynPol Main Loop - Step 2

**Policy Stack**



find_issue

find_repo

**Dynamic Policy Library**

create_group
find_subreddit
find_issue
search_customer
find_repo

**Main Agent** → outputs → Action

Relevant Policies

Page Operations

interacts

Policy Retrieval

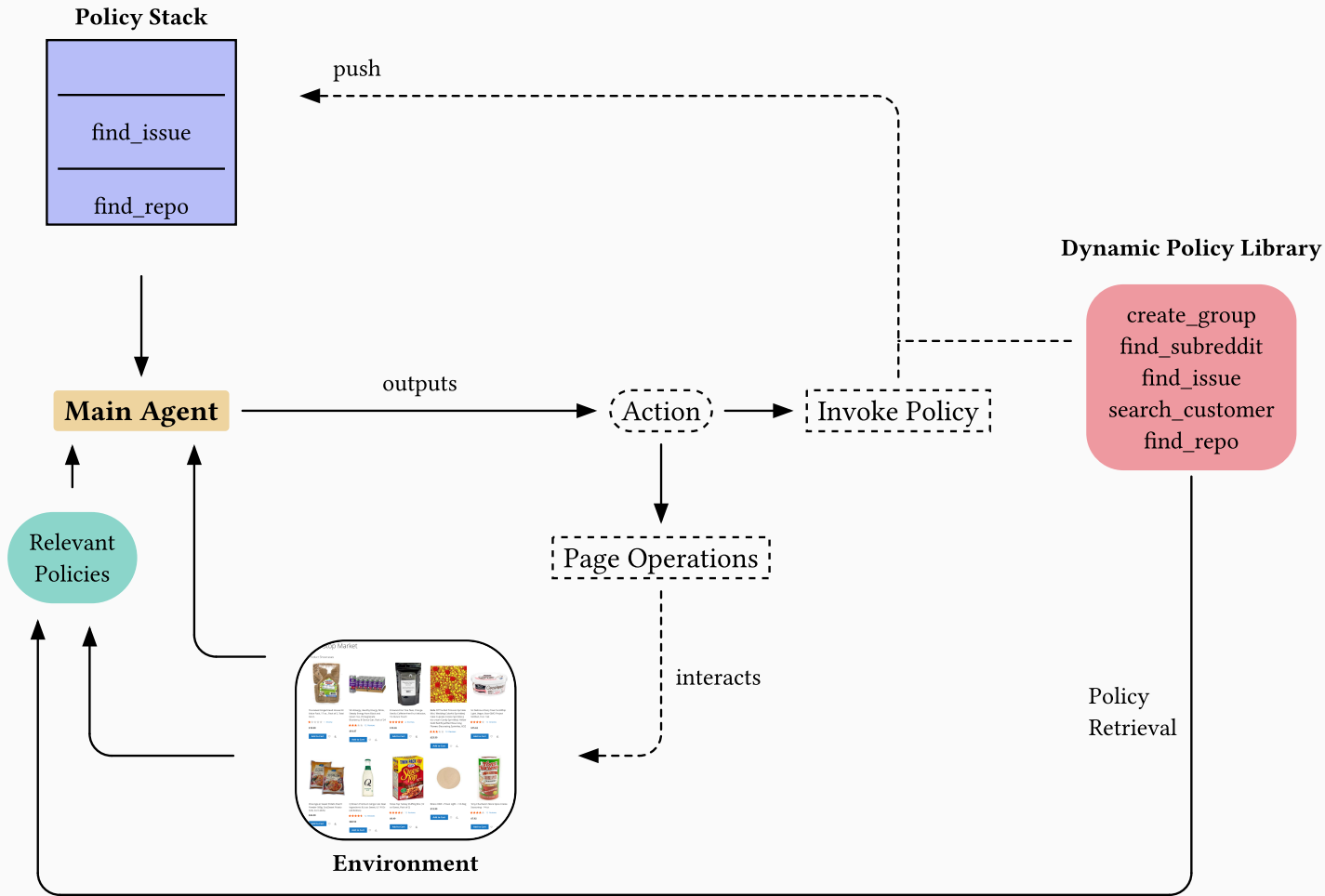**Environment**

First type of action : **Page Operation**

Atomic action that changes (or not) the environment. Examples :

- `click [380]`

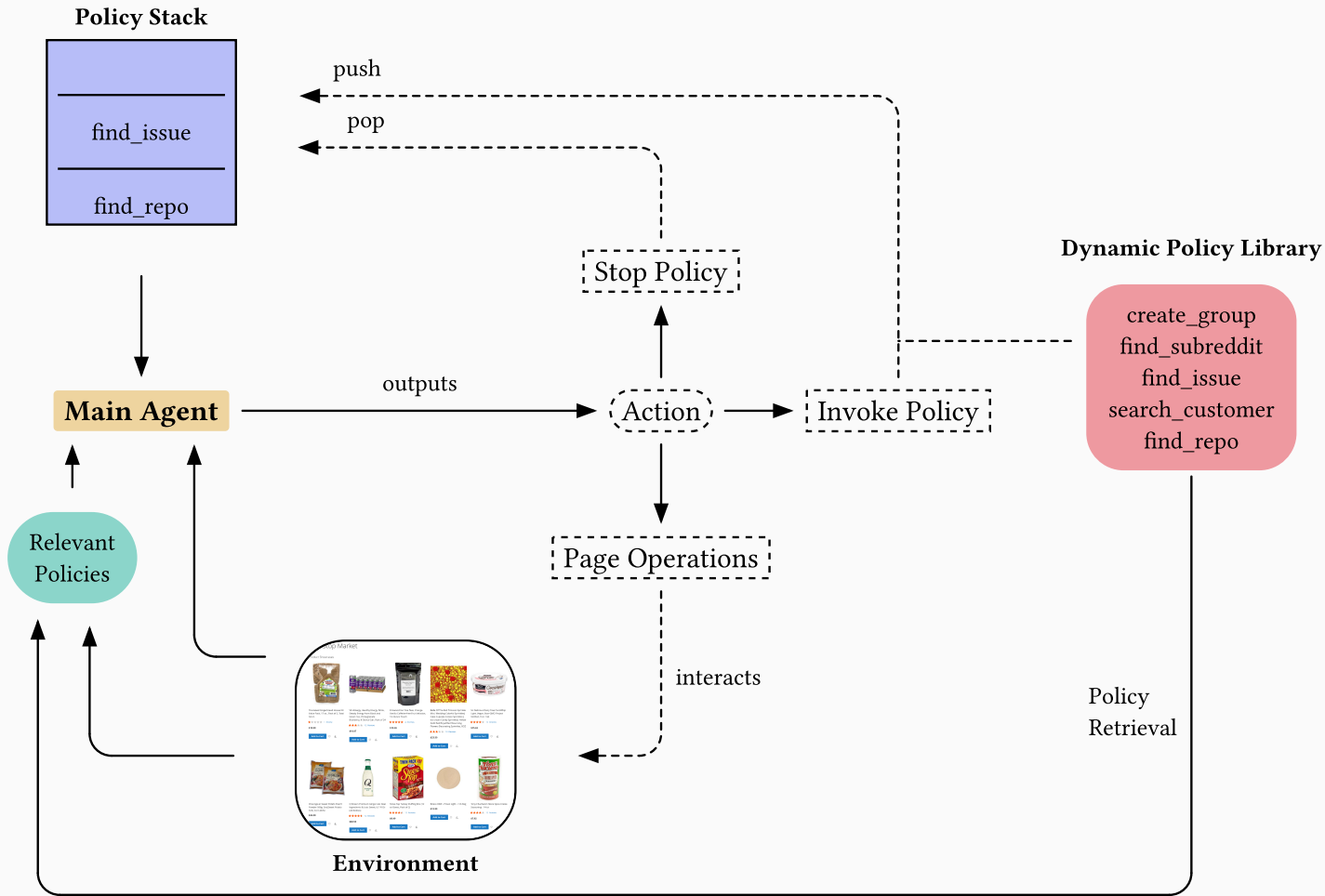- `type [67] [Carnegie Mellon]`

Figure 9: DynPol Main Loop - Step 3

Figure 10: DynPol Main Loop - Step 4

**Policy Stack**

find_issue

find_repo

push

**Dynamic Policy Library**

create_group
find_subreddit
find_issue
search_customer
find_repo

**Main Agent**

outputs

Action

Invoke Policy

Relevant Policies

Page Operations

interacts

Policy Retrieval

**Environment**

Second type of action : **Invoke a policy**

- Call a policy from **Relevant Policies**
- Add it on top of the stack

**Policy Stack**

find_issue

find_repo

push

pop

Stop Policy

**Dynamic Policy Library**

create_group
find_subreddit
find_issue
search_customer
find_repo

Main Agent

outputs

Action

Invoke Policy

Relevant Policies

Page Operations

interacts

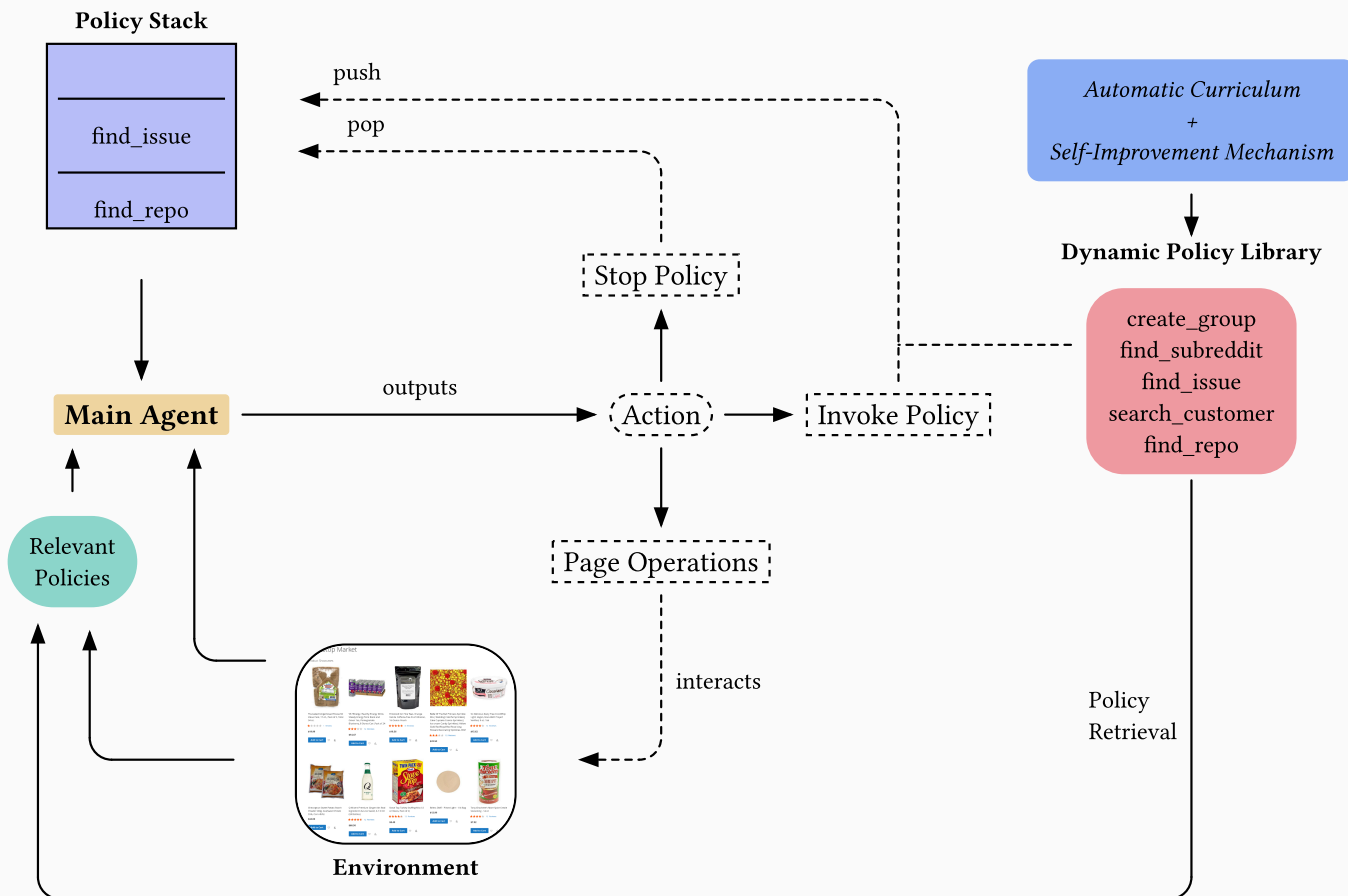Policy Retrieval

Environment

Figure 11: DynPol Main Loop - Step 5

Third type of action :
**Stop the currently active policy**

- **Pop** policy off the stack

- Return an answer

- If root policy :
**terminates** interaction

Figure 12: DynPol Main Loop - Step 5

- After finishing a task : **Self-Improvement Mechanism** : correct polices based on **interaction**
- Before starting a task : **Automatic Curriculum** : plan interaction and create new policies based on **environment**

Figure 13: Self-Improvement Mechanism

Trajectory + **Critique LLM** ⟹ Trajectory Breakdown

Trajectory Breakdown + Previous Instruction + **Rewriting LLM** ⟹ New instructions

Figure 14: Automatic Curriculum

Environment +
Relevant Policies +
**Automatic
Curriculum LLM** $\Rightarrow$
Plan

Plan $\Rightarrow$ Policies to
solve it (either new or
in **Relevant Policies**)

Improve DynPol through **experience**

More tasks encountered $\Rightarrow$ More policies in the **library** $\Rightarrow$ Better performance

Run DynPol on as many tasks as possible

Run DynPol on the same task multiple times:

- Better Refinement
- Track policies progress

# 3. Technical Details

## Webarena

- 5 types of websites, 812 tasks along with solutions
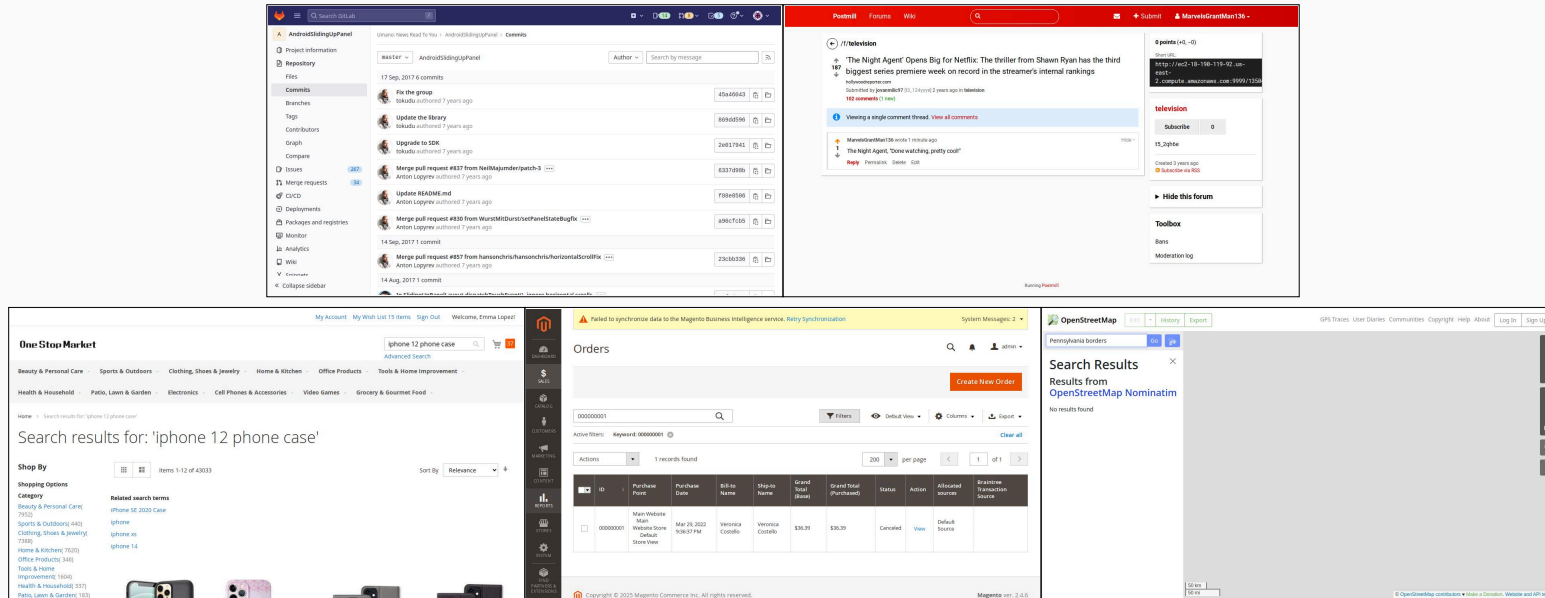- Self-Hosted, Closed, Controlled Environment (using AWS)



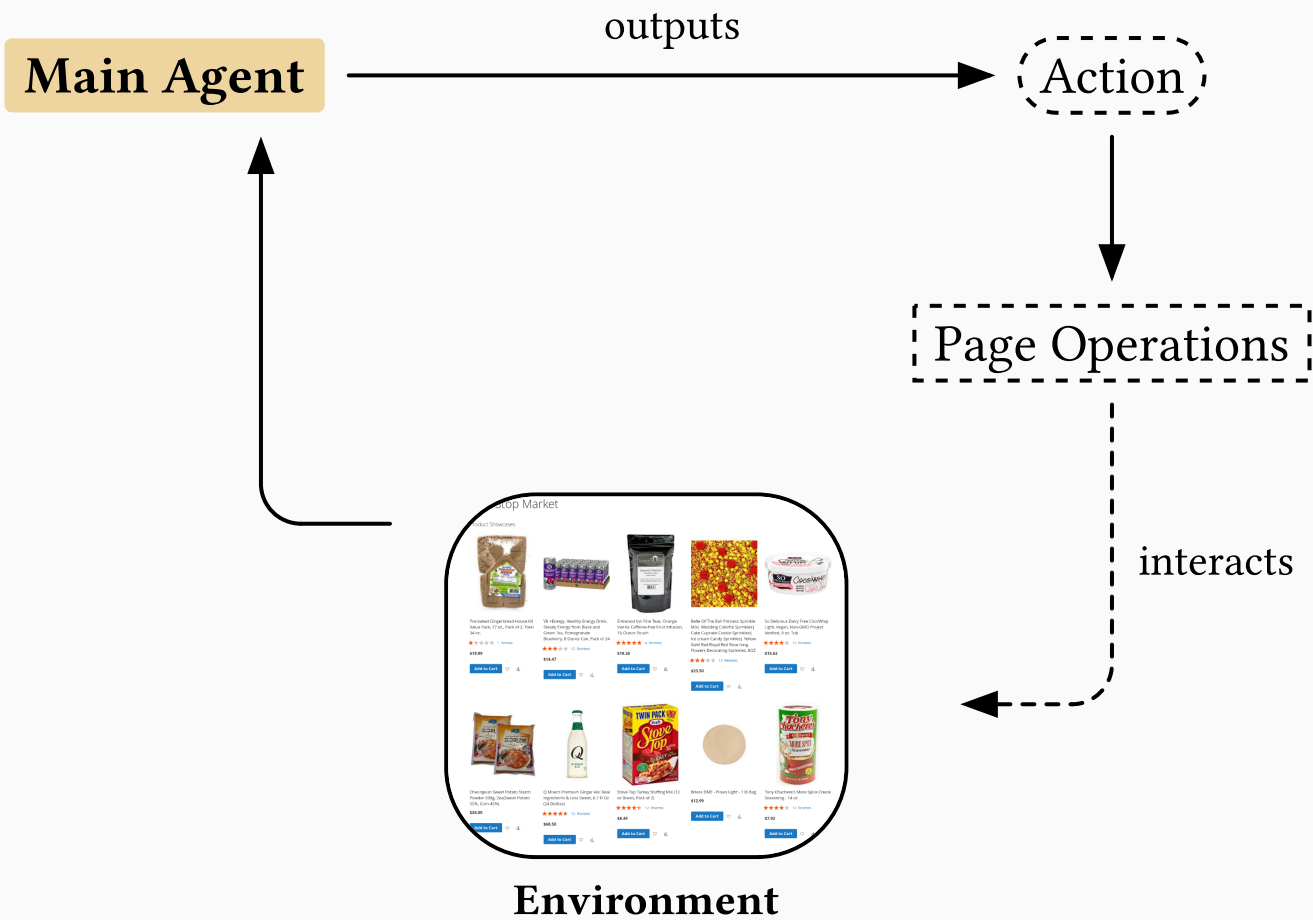Figure 15: From Left to Right, Top to Bottom :

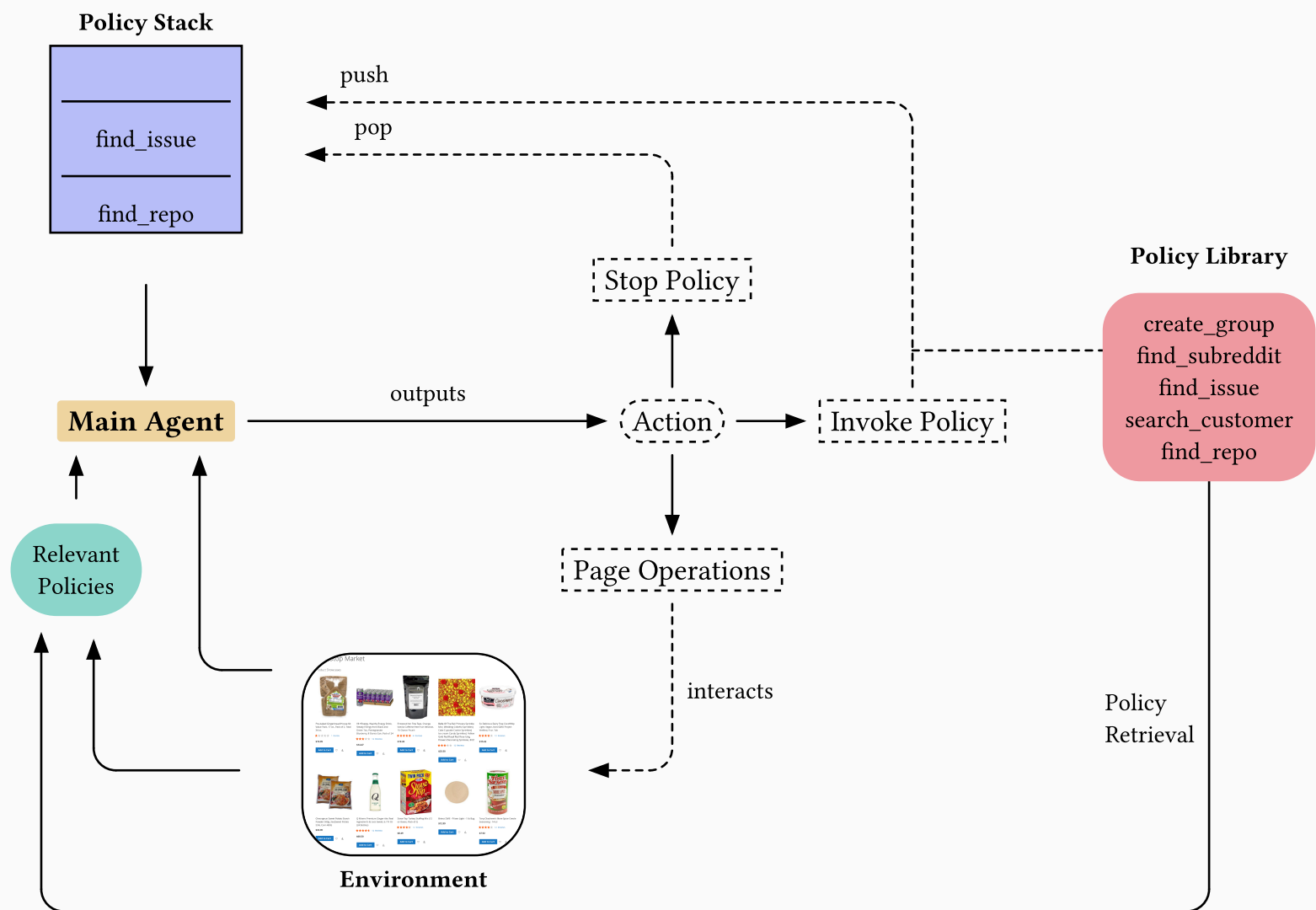GitLab, Reddit, Shopping, Shopping Admin, Maps (not working)

- **Main LLM** for experiments : `Llmama 3.3 Instruct 70B`

- LLM ran on the Marenostrum BSC with slurm

- Implementation made in python using the **Gymnasium** and **BrowserGym** libraries to interact with web agents

# 4. Baseline Comparison

All variants of DynPol are initialised with the **SteP policies**

Architectures used as baselines :

- DynPol without the automatic curriculum (i.e. only optimizing the SteP policies)

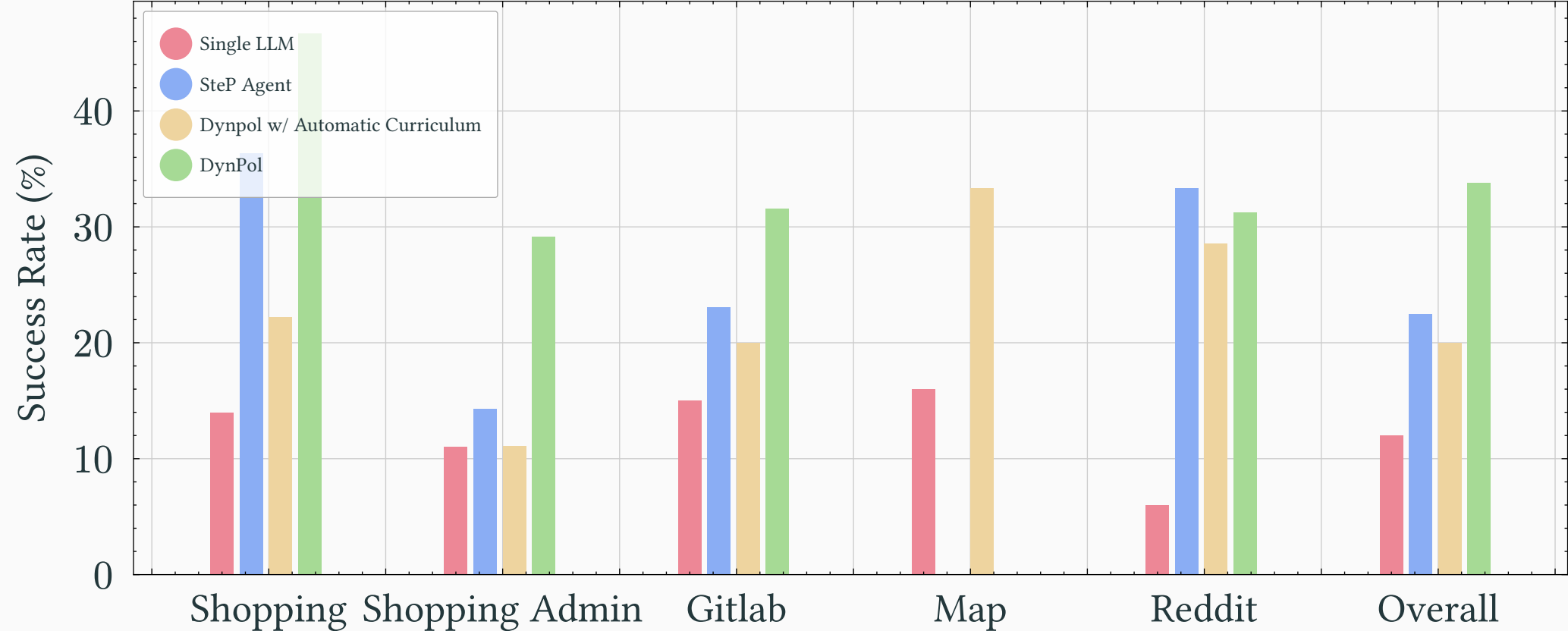- The DynPol architecture as presented before

| Tasks | Single LLM | SteP Agent | DynPol without Automatic Curriculum (3 iterations) | DynPol (3 iterations) |
|---|---|---|---|---|
| Shopping | 0.14 | 0.36 | 0.22 | **0.47** |
| Shopping Admin | 0.11 | 0.14 | 0.11 | **0.29** |
| Gitlab | 0.15 | 0.23 | 0.2 | **0.32** |
| Map | 0.16 | 0 | **0.33** | 0 |
| Reddit | 0.06 | 0.33 | 0.29 | **0.31** |
| Overall | 0.12 | 0.22 | 0.2 | **0.34** |

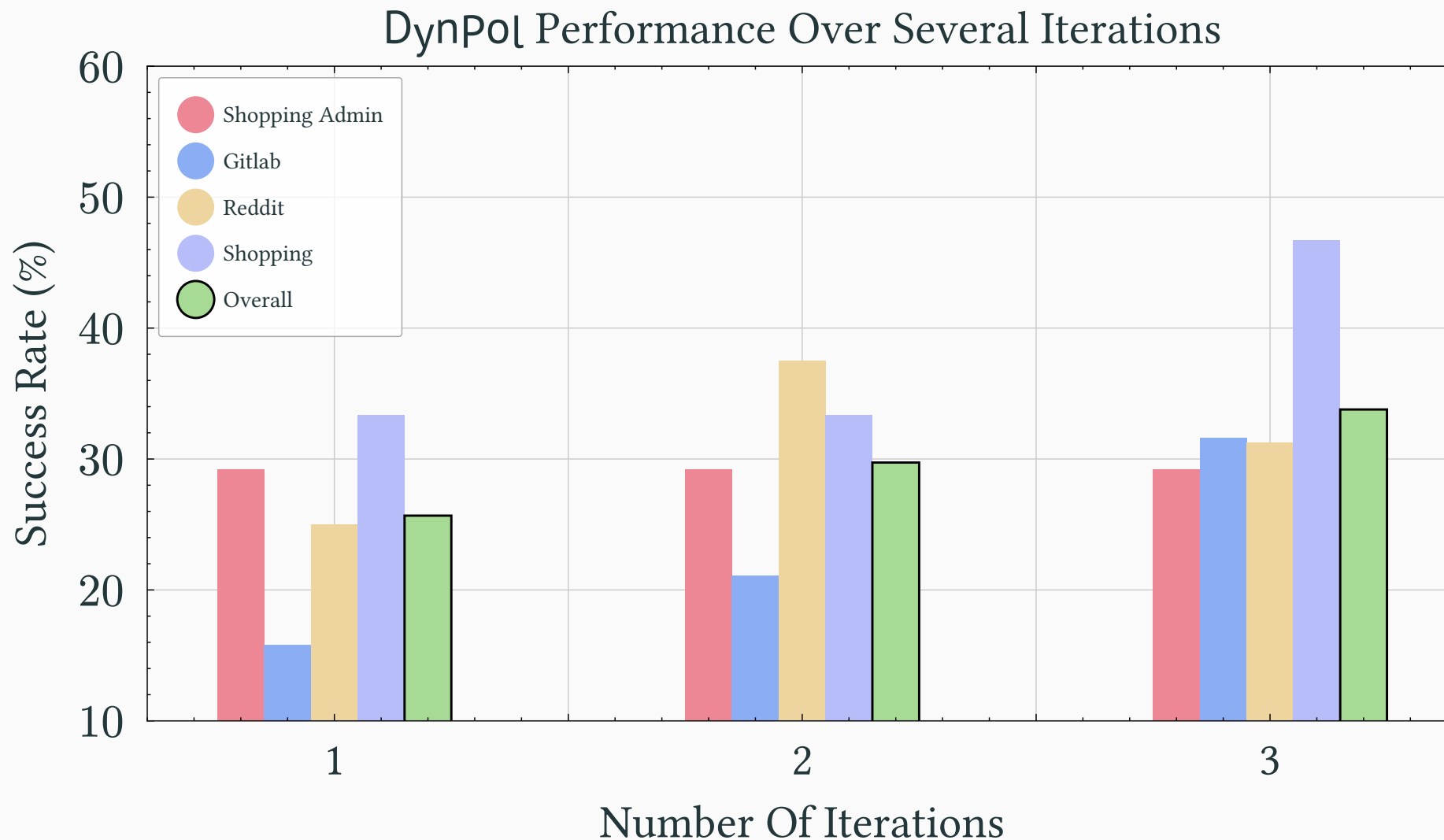Table 1: WebArena Success Rates for 50 (or 80) randomly sampled tasks for several architectures
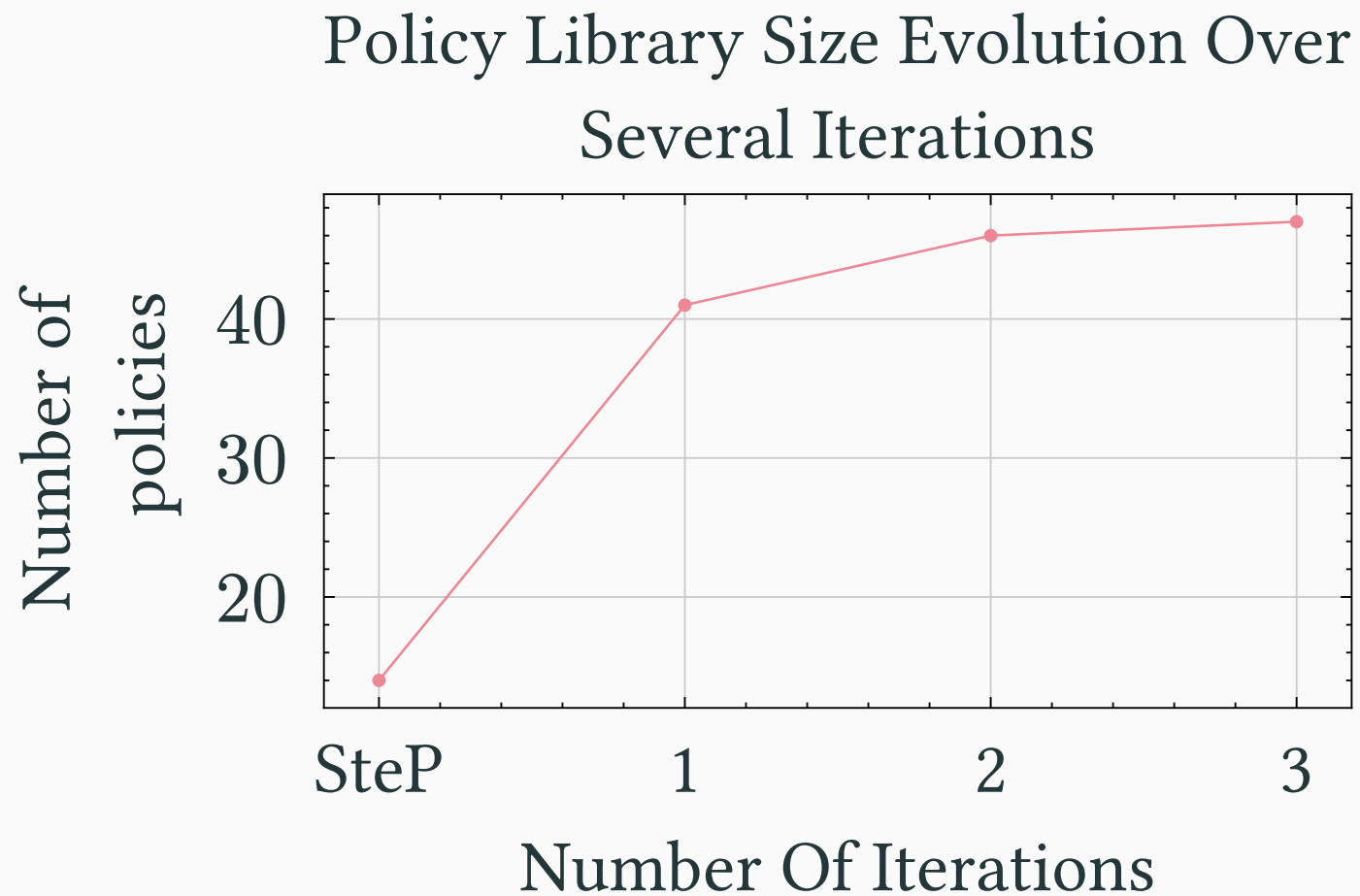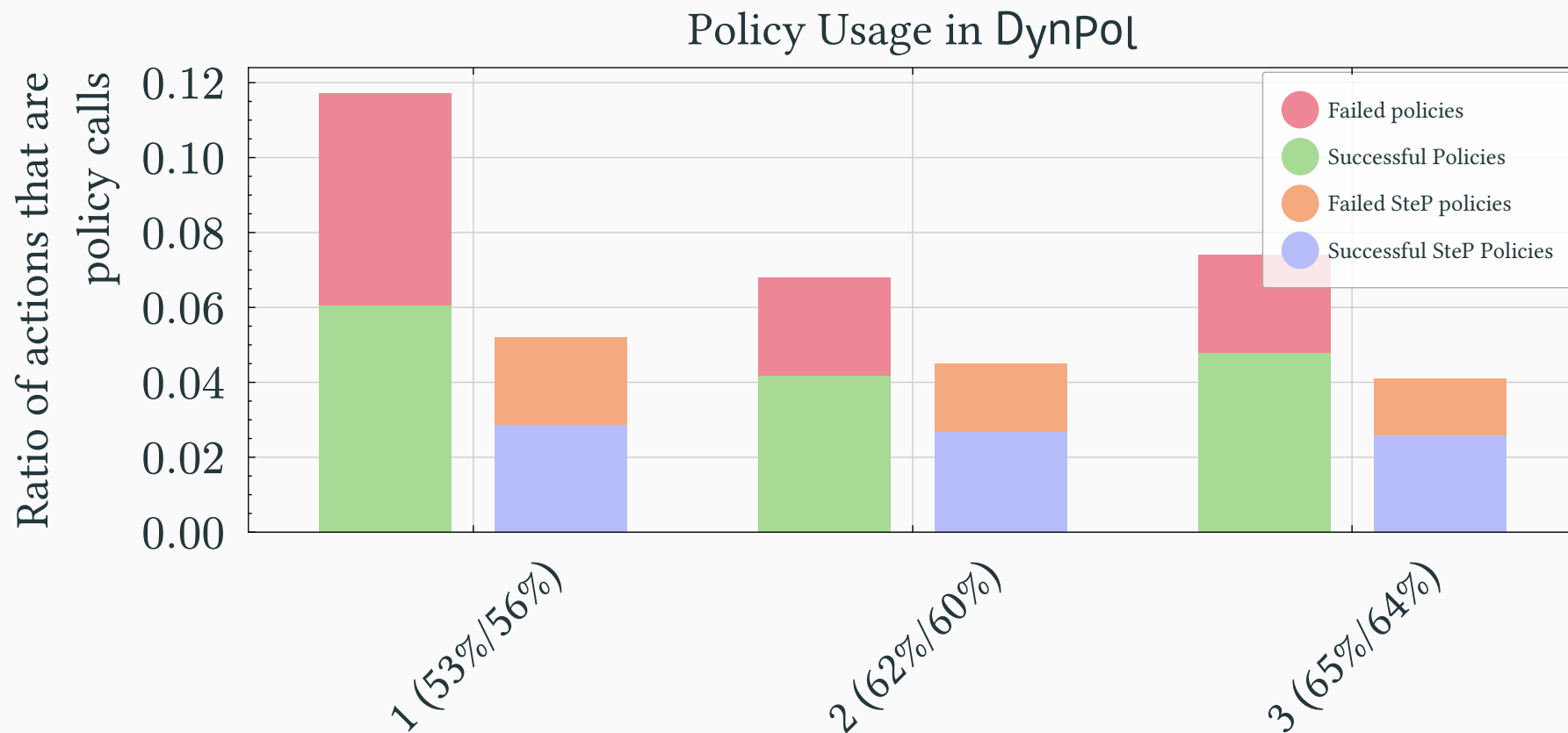
DynPol Performance Over Several Iterations

# 5. Analysis

Policy Library Size Evolution Over Several Iterations

Policy Usage in DynPol

# 6. Conclusion

# 6.1 Limitations

- Performance Gap

- Unreliability of LLM-driven prompt optimization

- Catastrophic forgetting in self-improvement

- Limitations of text-only approach

- Time-intensive execution

- More experiments

- Improve prompts

- Add visual components

# 7. Appendix

# Bibliography

[1] P. Sodhi, S. R. K. Branavan, Y. Artzi, and R. McDonald, "SteP: Stacked LLM Policies for Web Actions." [Online]. Available: https://arxiv.org/abs/2310.03720

[2] S. Marreed *et al.*, "Towards Enterprise-Ready Computer Using Generalist Agent." [Online]. Available: https://arxiv.org/abs/2503.01861

[3] S. Zhou *et al.*, "WebArena: A Realistic Web Environment for Building Autonomous Agents." [Online]. Available: https://arxiv.org/abs/2307.13854

[4] H. Su, R. Sun, J. Yoon, P. Yin, T. Yu, and S. Ö. Arık, "Learn-by-interact: A Data-Centric Framework for Self-Adaptive Agents in Realistic Environments." [Online]. Available: https://arxiv.org/abs/2501.10893

[5] K. Yang *et al.*, "AgentOccam: A Simple Yet Strong Baseline for LLM-Based Web Agents." [Online]. Available: https://arxiv.org/abs/2410.13825
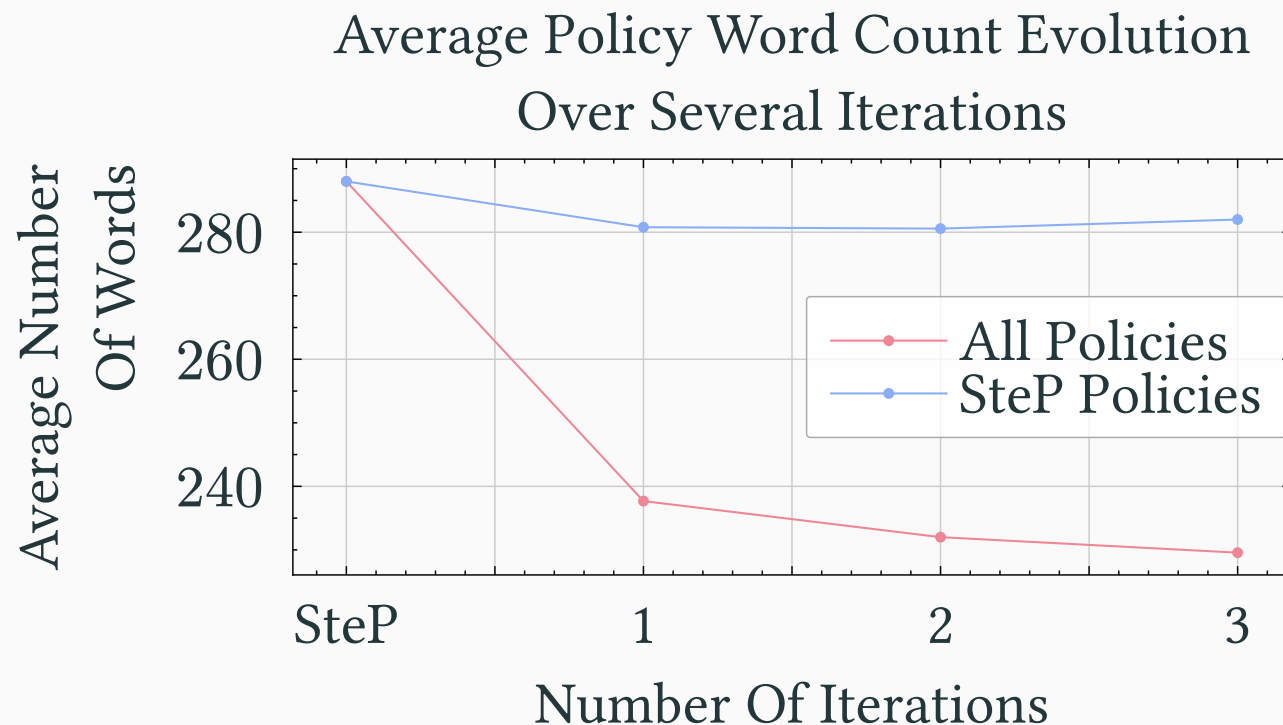
# Bibliography

[6]  S. Murty, H. Zhu, D. Bahdanau, and C. D. Manning, "NNetNav: Unsupervised Learning of Browser Agents Through Environment Interaction in the Wild." [Online]. Available: https://arxiv.org/abs/2410.02907

[7]  R. Zhang *et al.*, "Symbiotic Cooperation for Web Agents: Harnessing Complementary Strengths of Large and Small LLMs." [Online]. Available: https://arxiv.org/abs/2502.07942

[8]  G. Wang *et al.*, "Voyager: An Open-Ended Embodied Agent with Large Language Models." [Online]. Available: https://arxiv.org/abs/2305.16291

[9]  S. Murty, C. Manning, P. Shaw, M. Joshi, and K. Lee, "BAGEL: Bootstrapping Agents by Guiding Exploration with Language." [Online]. Available: https://arxiv.org/abs/2403.08140

[10]  T. Khot *et al.*, "Decomposed Prompting: A Modular Approach for Solving Complex Tasks." [Online]. Available: https://arxiv.org/abs/2210.02406

# Bibliography

[11] M. Ghasemi and D. Ebrahimi, "Introduction to Reinforcement Learning." [Online]. Available: https://arxiv.org/abs/2408.07712

[12] M. Yuksekgonul *et al.*, "TextGrad: Automatic "Differentiation" via Text." [Online]. Available: https://arxiv.org/abs/2406.07496

[13] Y. Lu, J. Yang, Y. Shen, and A. Awadallah, "OmniParser for Pure Vision Based GUI Agent." [Online]. Available: https://arxiv.org/abs/2408.00203

[14] S. Yao *et al.*, "ReAct: Synergizing Reasoning and Acting in Language Models." [Online]. Available: https://arxiv.org/abs/2210.03629

[15] J. Shen *et al.*, "ScribeAgent: Towards Specialized Web Agents Using Production-Scale Workflow Data." [Online]. Available: https://arxiv.org/abs/2411.15004

## Average Policy Word Count Evolution Over Several Iterations



SteP policies word count < DynPol - generated policies word count

SteP policies' instructions contain toy examples

# 7.3 Pseudo-Code of DynPol Main Loop

```python
def get_action(observation, objective, is_first_step):
    if is_first_step:
        # Initialize the policy stack with a site-dependant policy
        policy_stack = Stack.init(root_policy)
    # Automatic Curriculum, seeing if new policies need to be added
    # to the policy library
    relevant_policies = policy_library.retrieve(objective)
    # LLM Call to choose which policies to use
    chosen_policies = automatic_curriculum(objective, observation,
    relevant_policies)
    for policy in chosen_policies:
        if not (policy in policy_library):
            policy_library.add(policy)

```

```
13      current_policy = policy_stack.top()
14      relevant_policies = policy_library.retrieve(objective)
15      # LLM call to determine which action to take
16      action = get_action(objective, observation, current_policy,
        relevant_policies)
17
18      match action.type:
19          case PAGE_OPERATION:
20              env.interact(action)
21          case POLICY_CALL:
22              policy_stack.push(action.policy)
23          case STOP:
24              finished_policy = policy_stack.pop()
25              # Self-Improvement Mechanism
```

```
26    critique = get_critique(finished_policy.trajectory,
      objective)
27    policy_library.update_usage(critique.is_success)
28    if
      rewrinting_needed(policy_library.get_usage(finished_policy)):
29        new_guidance = rewrite(finished_policy.guidance,
          critique, objective, finished_policy.trajectory)
30        policy_library.update(finished_policy, new_guidance)
```
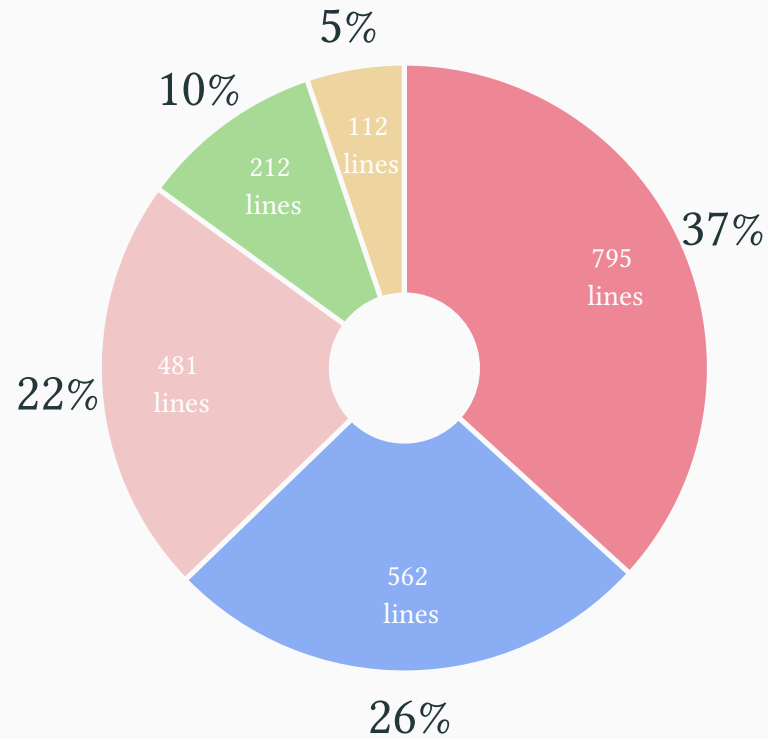
# 7.4 Sample Policy Prompt

Based on the plan and explanation, the guidance for the web agent to solve the task is:
* To find the latest post related to a specific topic, start by examining the search results page and looking for the most recent timestamp.
* Use the website's search function or filtering capabilities to narrow down the search and highlight the latest post. Look for options to sort posts by timestamp in descending order (newest first).
* If sorting or filtering options are available, use them to highlight the latest post.
* Clearly identify the latest post and its author by looking for the username associated with the post.
* Verify that the highlighted post is indeed the latest one related to the topic by checking the timestamp and ensuring that the post is relevant to the search query.
* Provide clear evidence in the observation that the task has been completed by highlighting the latest post and its author.

Examples:
* If the website has a search box, type in the keyword to find posts related to the topic, and then look for the most recent timestamp to identify the latest post.
* If the website has filtering options, select the option to filter posts by keyword or topic, and then select the topic to find related posts. Look for options to sort posts by timestamp in descending order.
* If the website displays user posts in a list, look for the most recent timestamp to identify the latest post related to the topic, and then extract the author's username.
* If the website displays user posts in a grid, look for the post with the most recent timestamp, extract the author's username, and verify that the post is relevant to the search query.

Figure 16: Repartition of the lines inside the codebase

**March**

- Literature for the initial topic
- Switching topic
- Literature for the actual topic
- Initial DynPol draft

**April**

- Implementing DynPol
- Implementing tools (logging, interacting with web env,...)

**May**

- Refining DynPol (Automatic Curriculum details, Critique, ...)
- Implementing baseline
- Starting to setup MareNostrum access

## June

- Experiments
- Refining DynPol
- Set up our LLM on MareNostrum
- Fixing the WebArena AWS instance

## July

- More Experiments
- Prompt Engineering on DynPol
- Fixing MareNostrum problems
- Starting report

## August

- Final Experiments
- Finishing report