

Agent-Based SRE: Automated Diagnosis and Mitigation in K8s

Marcello Martini

Problem

Modern cloud-native infrastructures (e.g., Kubernetes) are complex and prone to faults.

Site Reliability Engineers (SREs) face challenges in rapidly diagnosing and mitigating incidents in distributed systems.



There is a need for intelligent systems that can support SREs with actionable insights and automated mitigation plans

Goal

Develop a prototype of an autonomous agent that observes the overall architecture, identifies issues, and proposes (or initiates) mitigation steps in Kubernetes environments.



State of the Art

AIOpsLab - Microsoft (<https://doi.org/10.48550/arXiv.2501.06706>)

AIOpsLab presents a benchmark suite for evaluating AI agents in autonomous cloud management, comparing LLM-based and traditional AIOps solutions using fault-injected microservices scenarios.

ITBench - IBM (<https://doi.org/10.48550/arXiv.2502.05352>)

ITBench introduces an open-source framework for benchmarking AI agents on diverse, real-world IT automation tasks, supporting multiple IT personas and offering comprehensive observability and alerting.

MonitorAssistant - Microsoft (<https://doi.org/10.1145/3663529.3663826>)

MonitorAssistant proposes an end-to-end system that leverages large language models to simplify cloud service monitoring, providing automated configuration recommendations, practical anomaly detection, and interactive troubleshooting guidance.

Tools



Microsoft AIOpsLab

Microsoft AIOpsLab is used to set up the testbed (i.e., the distributed infrastructure) and inject faults.



LangChain & LangGraph

LangChain and LangGraph are used to create workflows and agents.



LLMs

Large Language Models are used as engines for parsing and processing metrics.



MCP servers

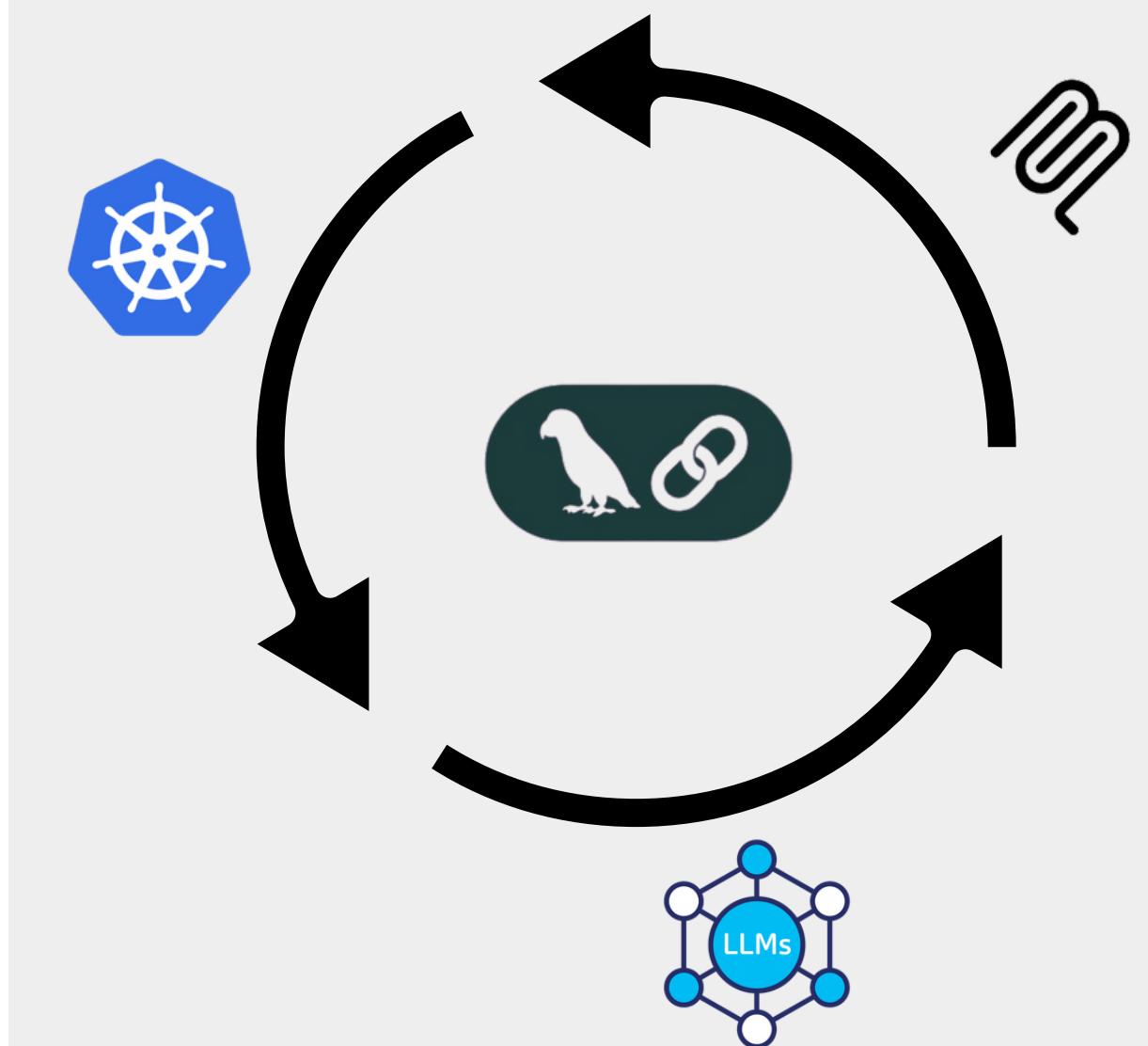
Model Context Protocol servers are used to interact with Kubernetes and other external resources.

Solution Architecture

Agent-based system with modular design.

Langchain agent connects to multiple MCPs:

- Kubernetes MCP: Cluster state and events
- Prometheus MCP: Metrics and alerts
- RAG MCP: Knowledge retrieval
- (expandable)



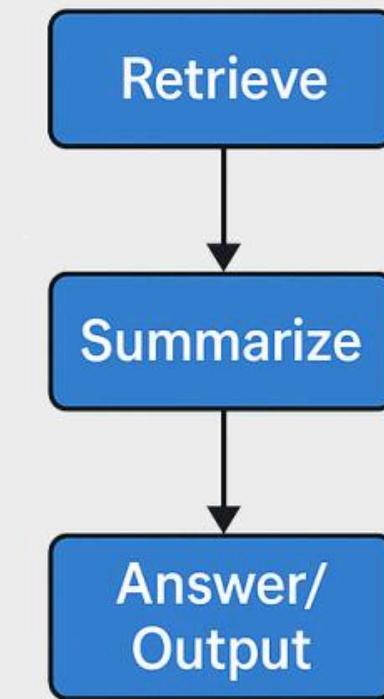
Requirements

- Agent must autonomously detect and diagnose issues.
- Provide clear explanations for each step taken.
- Suggest and/or execute mitigation actions.
- Evaluation framework to assess agent performance.

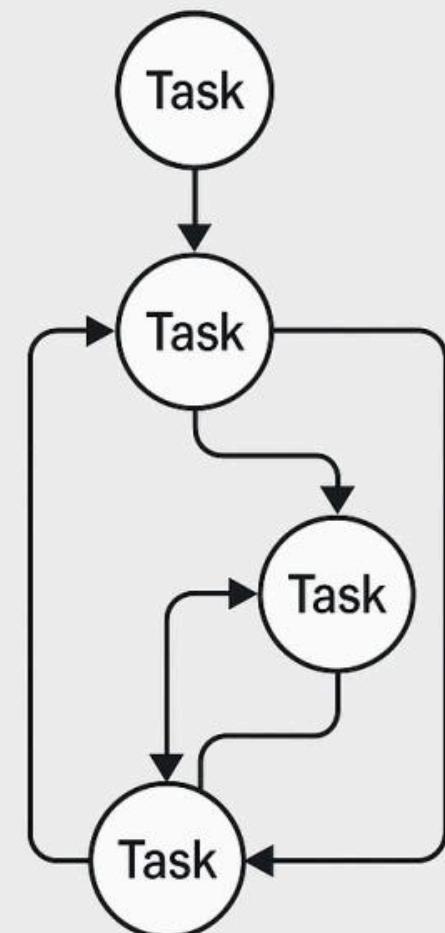
LangChain and LangGraph

LangChain is a framework that simplifies the development of powerful applications using large language models by connecting them with data sources, tools, and workflows.

LangGraph builds on LangChain, enabling developers to create dynamic, multi-step conversational agents with flexible, **graph-based workflows**.



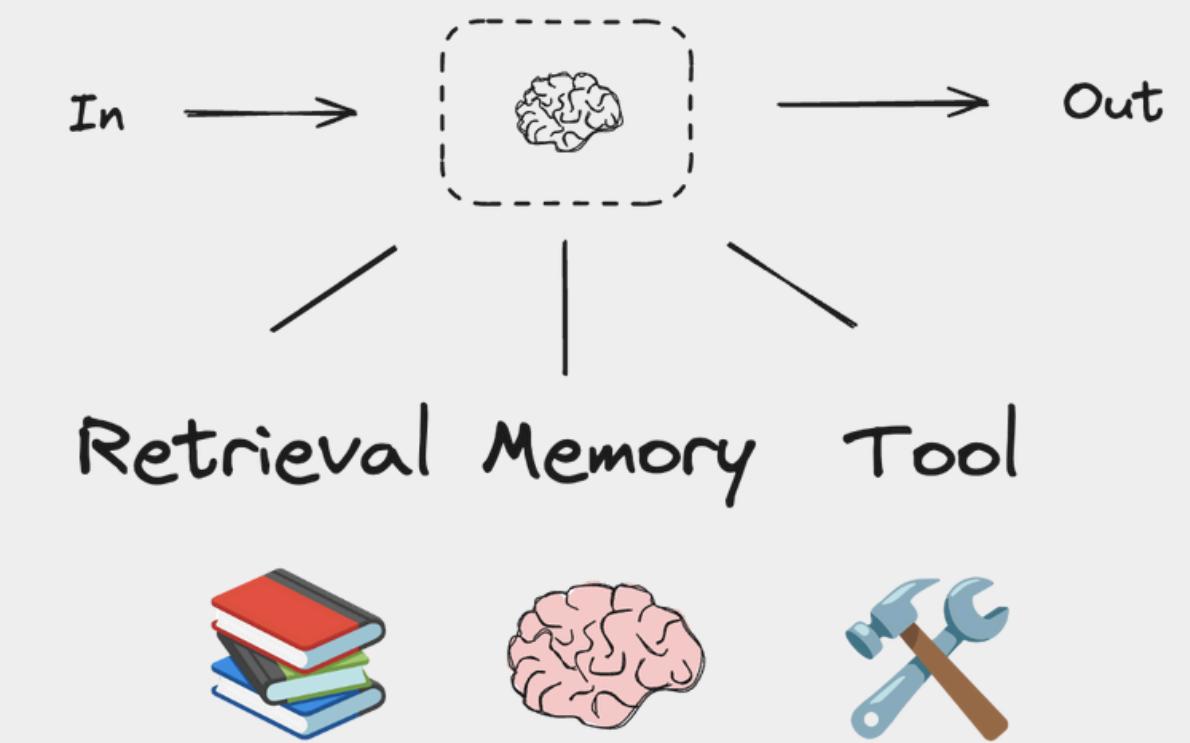
LangChain



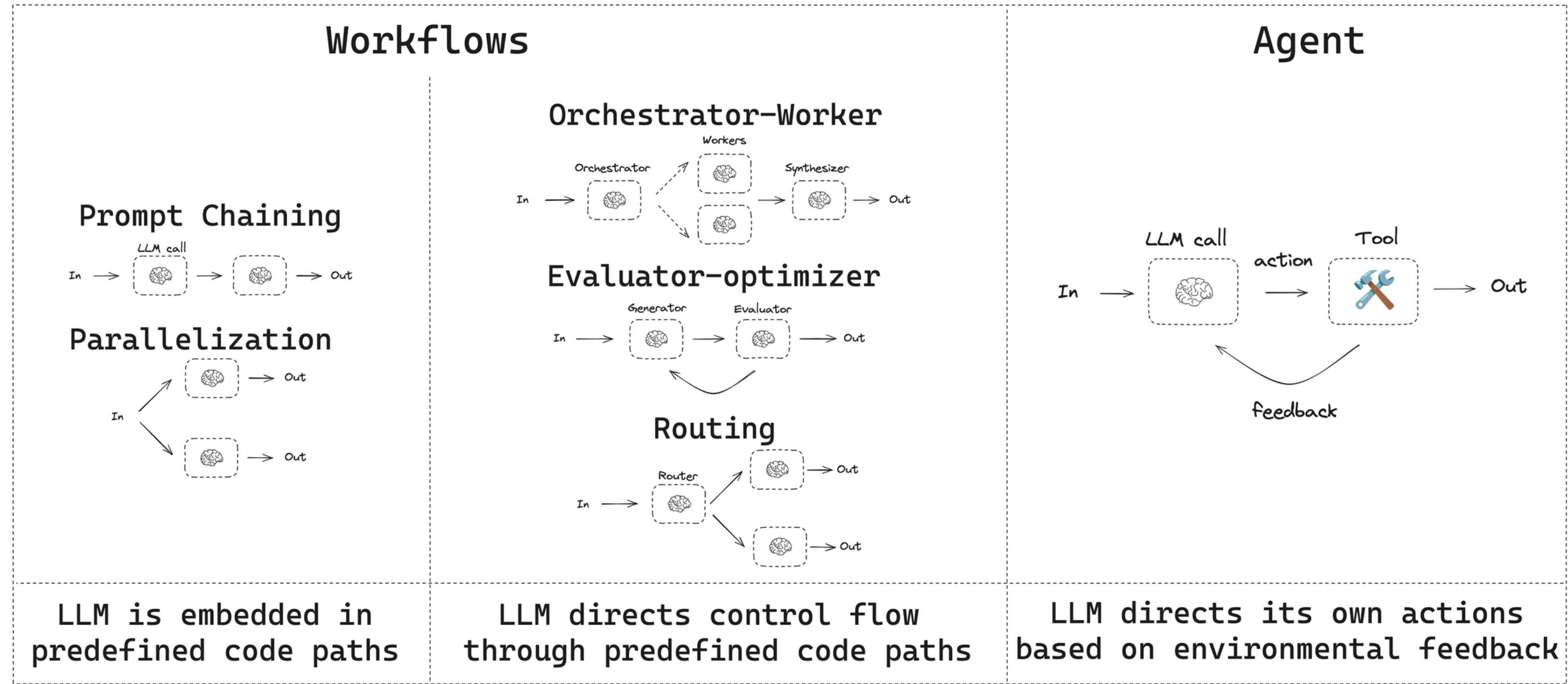
LangGraph

LangChain and LangGraph

- **Modular LLM Workflows:** Build complex applications by chaining together language models and tools.
- **Dynamic Conversational Agents:** Create agents capable of multi-step reasoning.
- **Rich Ecosystem Integrations:** Connect easily with databases, APIs, document loaders etc.
- **State Management:** Maintain conversation history and manage state across interactions.



LangChain and LangGraph



MCP servers

A GitHub repository card for the Flux159/mcp-server-kubernetes repository. It features a blue octagonal icon with a white steering wheel. The repository name is "Flux159/mcp-server-kubernetes". Below it, the description "MCP Server for kubernetes management commands" is shown. The stats are: 23 contributors, 7 used by, 903 stars, and 140 forks. A large blue letter "S" is positioned to the right of the card.

A GitHub repository card for the chroma-core/chroma-mcp repository. It features a yellow, red, and blue overlapping circle icon. The repository name is "chroma-core/chroma-mcp". Below it, the description "A Model Context Protocol (MCP) server implementation that provides database capabilities for Chroma" is shown. The stats are: 4 contributors, 4 issues, 241 stars, and 44 forks. A large blue letter "S" is positioned to the left of the card. Below the card, the text "Chroma DB is a vector database designed for the efficient storage and retrieval of vector embeddings (for RAG purposes)" is displayed.

A GitHub repository card for the pab1it0/prometheus-mcp-server repository. It features a red flame icon. The repository name is "pab1it0/prometheus-mcp-server". Below it, the description "A Model Context Protocol (MCP) server that enables AI assistants to query and analyze Prometheus metrics through standardized interfaces" is shown. The stats are: 6 contributors, 1 issue, 157 stars, and 32 forks. A small profile picture of a man is shown next to the card.

A GitHub repository card for the mendableai/firecrawl repository. It features a fire icon. The repository name is "mendableai/firecrawl". Below it, the description "Turn entire websites into LLM-ready markdown or structured data. Scrape, crawl and extract with a single API." is shown. The stats are: 98 contributors, 144 used by, 39 discussions, 43k stars, and 4k forks. A small profile picture of a person is shown next to the card. Below the card, the text "Firecrawl is an API that crawls websites to extract clean, structured content for AI use." is displayed.

Large Language Models



OpenAI provides access to the advanced **GPT-4.1** model, offering 250,000 free tokens daily



Google AI Studio enables experimentation with **Gemini models** and includes a free tier for select versions (flash)



Ollama

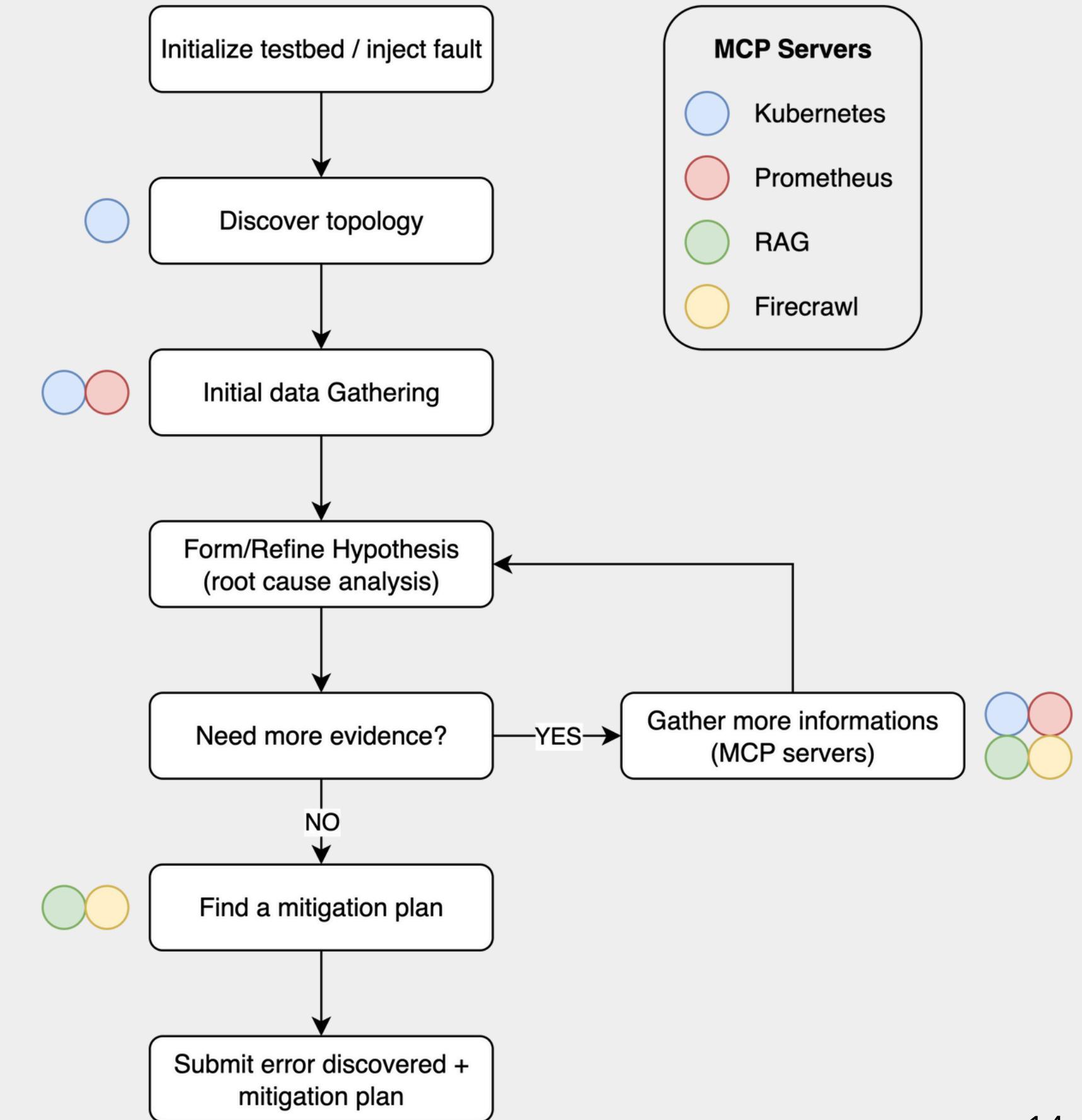
Ollama runs models **locally**, but not viable here due to hardware constraints.

Workflow

How to design the workflow?

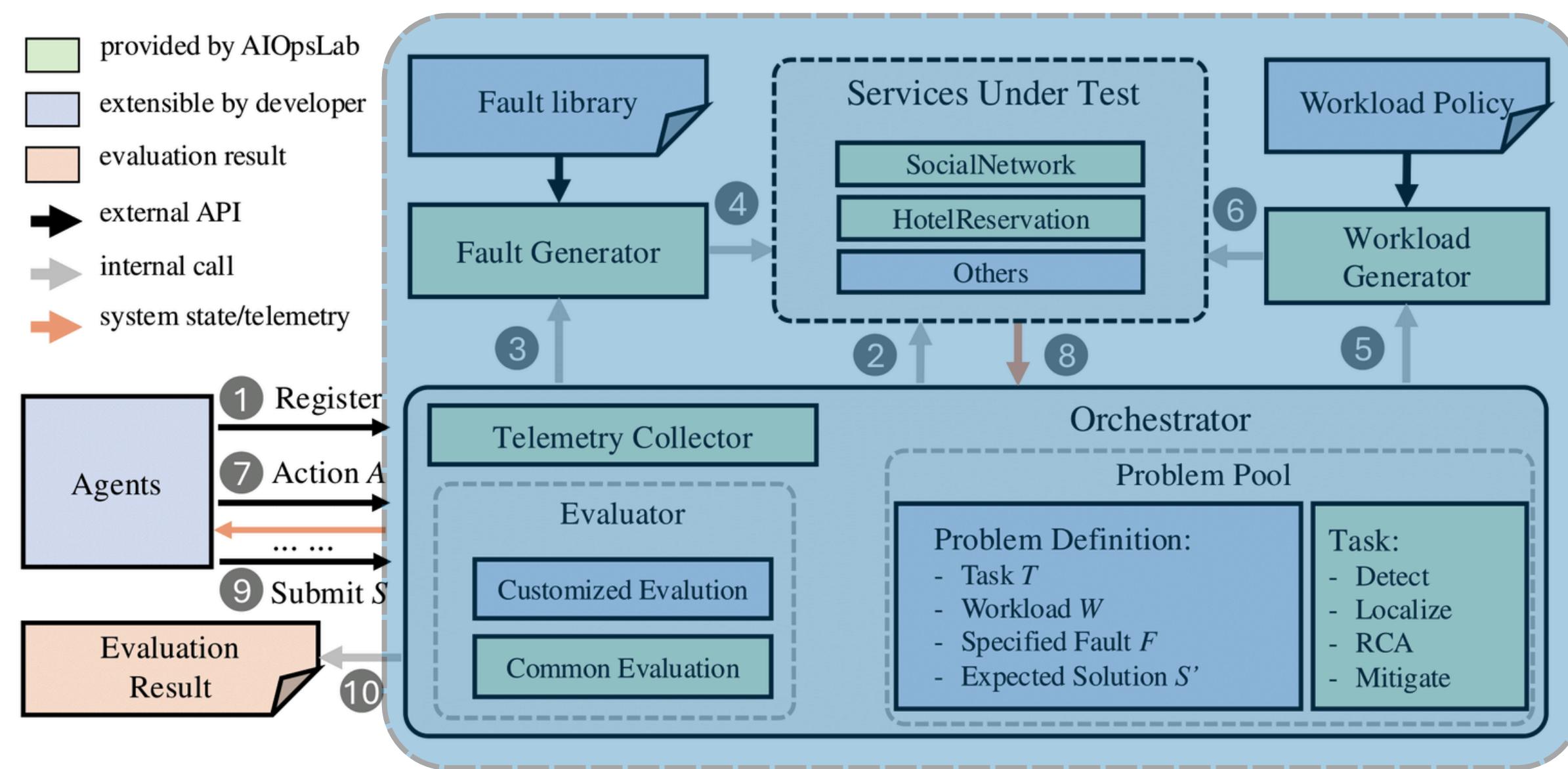
Need to find a trade off between effectiveness and simplicity.

Critical point: how to implement the feedback loop without human in the loop



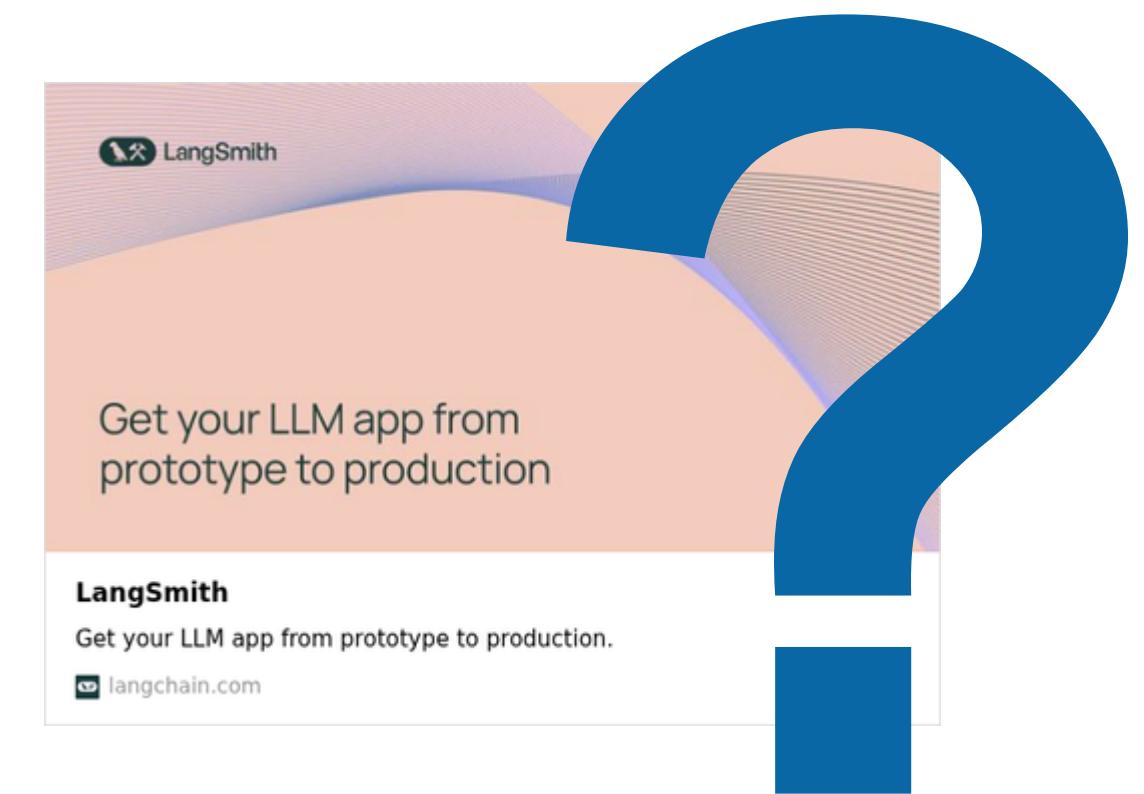
Testbed & Benchmarking

Microsoft AIOpsLab used to generate realistic infrastructure and inject faults (e.g., misconfigurations, DoS)



Performance evaluation

- Root cause detected (correctness)
- Total tokens consumed to solution
- Total steps executed to reach solution
- Elapsed time to solution
- Estimated cost per simulation



Open questions

- How to scale agent reasoning for large-scale clusters?
- How to optimize retrieval and RAG integration?
- Can the agent autonomously adapt to novel fault types?
- Integration with vector DBs for richer context



Potential thesis topics

- Can small models (SLMs) perform as well as larger ones if the workflow is organized differently?
- Are LLMs sufficient for understanding and detecting faults in microservices architectures?
- Can model distillation be effective for specific tasks in IT operations?



Thank you