# AgentQuest: A Modular Benchmark Framework to Measure Progress and Improve LLM Agents

*Luca Gioacchini, Giuseppe Siracusano, Davide Sanvito, Kiril Gashteovski, David Friede, Roberto Bifulco, Carolin Lawrence*

## Benchmarking Generative AI Agents

- o **Generative AI Agents**: software systems leveraging LLMs to perform complex multi-steps tasks
  - ▪ Key concepts: *Environment, State, Observation, Action*
- o Benchmarking is essential for evaluation and guiding improvements
- o **Issues in current benchmarks:**
  - ▪ Limited coverage of different types of benchmarks
  - ▪ Existing metrics focus on final success:
    - ▪ **Success Rate (SR)** and **Time to Success (Steps)**
  - ▪ Coupling of benchmarks and agent architectures

> ❌ Closed-box vs open-ended tasks with tools usage
> ❌ Limited insights into *intermediate* success/failure
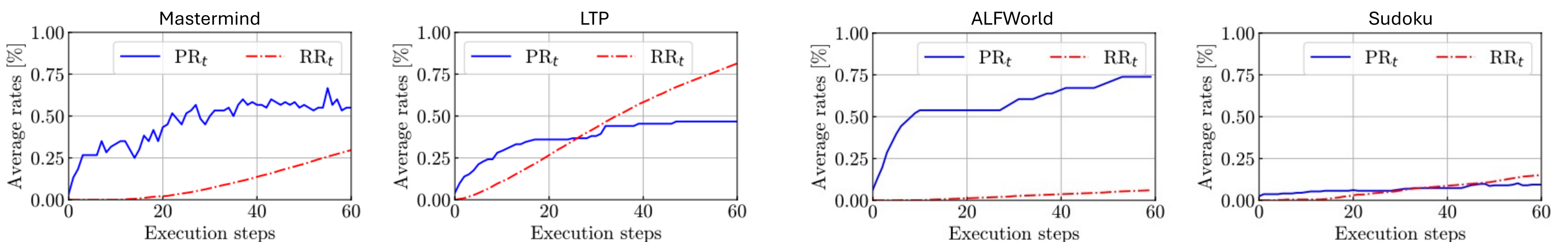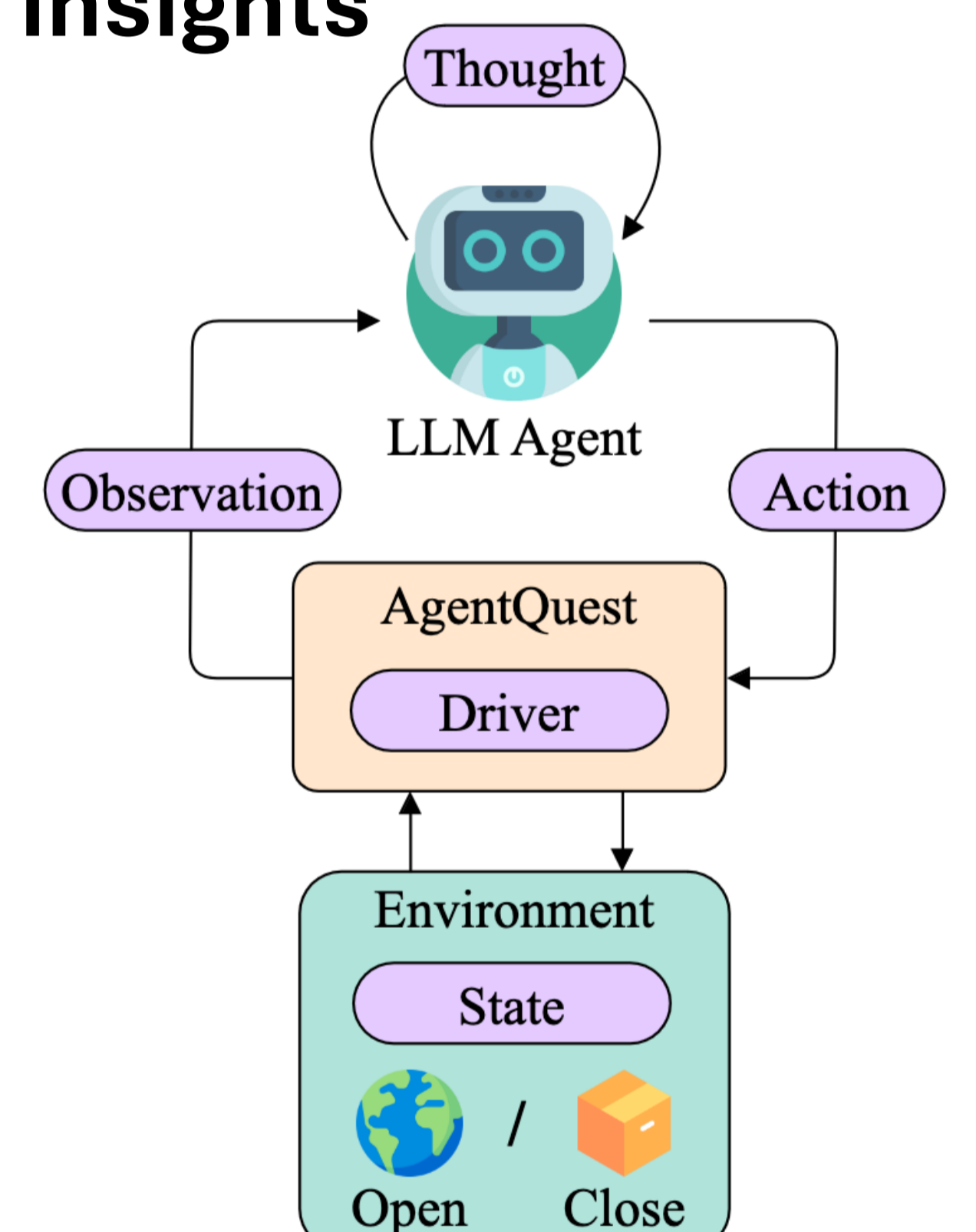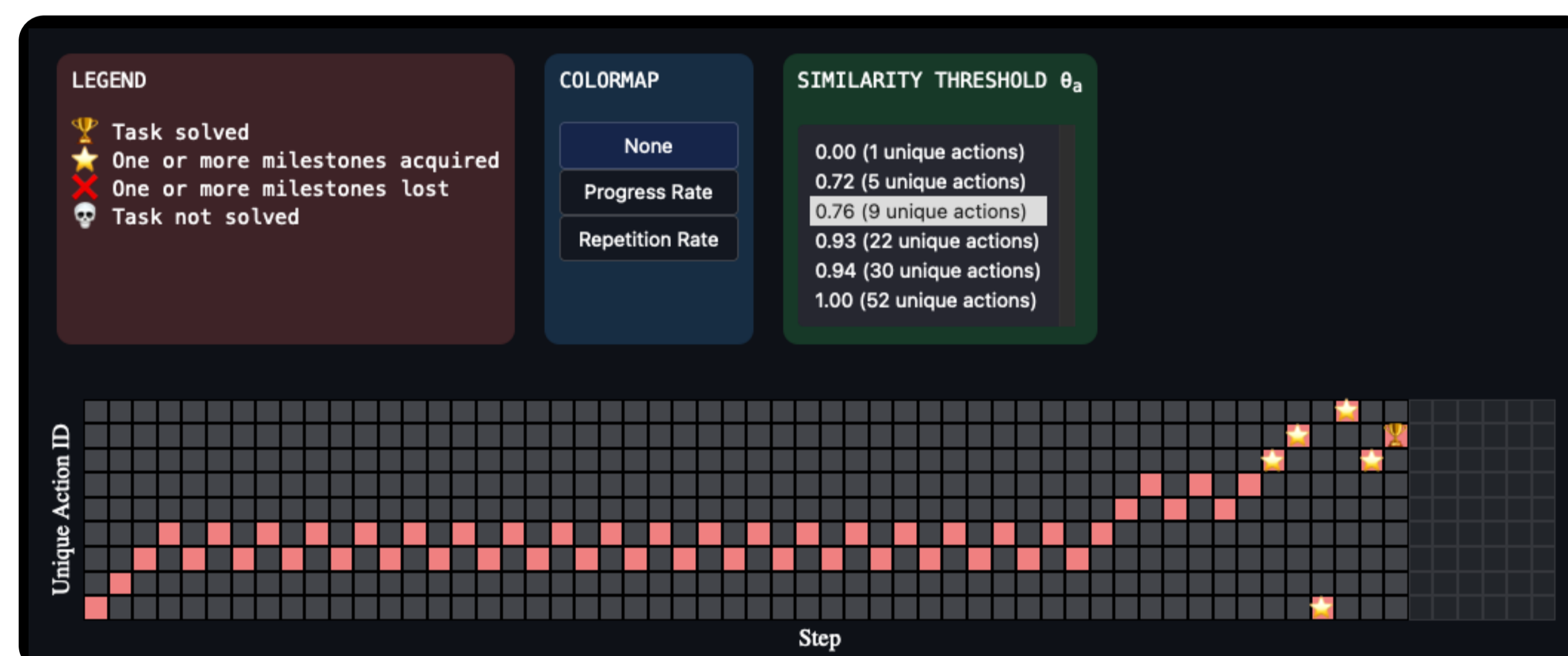> ❌ Limited re-usability and extensibility

## AgentQuest

- o Moudular framework to support multiple diverse benchmarks and agent architectures
- o Single unified Python interface
- o Two new cross-benchmarks metrics:
  - ▪ **Progress Rate (PR)** and **Repetition Rate (RR)**
- o Enables the **investigation of agent behaviour** to support the improvement of its architecture
- o (Initial) support for 4 benchmarks:
  - ▪ Evaluation of an agent based on LangChain Chat Model powered by GPT-4
  - ▪ *Improved results** after modifying the agent architecture **based on AgentQuest's insights**
- o Much more in the paper and in the GitHub repository!
  - ▪ Additional agent architectures and open-ended benchmark

> ✅ Agent / Environment decoupling
> ✅ Increased progress observability
> ✅ Simplified interoperability across benchmarks and agent architectures and improved extensibility



AgentQuest high-level architecture.

| | Existing Metrics | | AgentQuest | |
|---|---|---|---|---|
| | SR | Steps | $PR_{60}$ | $RR_{60}$ |
| Mastermind | 0.47 | 41.87 | 0.62 | 0.32 |
| LTP | 0.20 | 52.00 | 0.46 | 0.81 |
| ALFWorld | 0.86 | 21.00 | 0.74 | 0.06 |
| Sudoku | 0.00 | 59.67 | 0.08 | 0.22 |
| Mastermind* | 0.60 +13%pts | 39.73 | 0.73 +11%pts | 0.00 |
| ALFWorld* | 0.93 +7%pts | 25.86 | 0.80† +6%pts | 0.07† |

†Metrics with extended runtime up to 120 steps, i.e. $PR_{120}$ and $RR_{120}$.



The **AgentQuest GUI** provides a convenient visualization to track the progress over time and perform a step-by-step investigation of the actions performed by the agent.



AgentQuest enables to monitor progress and action repetition to get relevant insights to improve agent architecture: e.g. adding a memory component (for Mastermind) or extending the benchmark runtime (for ALFWorld).

## AgentQuest is open-source!

- 🌐 https://github.com/nec-research/agentquest
- 📖 Documentation for existing benchmarks
- 💡 Examples based on OpenAI and LangChain APIs
- 🚀 HOWTOs for new benchmarks, drivers and metrics

VIDEO

CODE

NEC · Politecnico di Torino · NAACL 2024