

Novelty Assessment Report

Paper: Gaia2: Benchmarking LLM Agents on Dynamic and Asynchronous Environments

PDF URL: <https://openreview.net/pdf?id=9gw03JpKK4>

Venue: ICLR 2026 Conference Submission

Year: 2026

Report Generated: 2025-12-27

Abstract

We introduce **Gaia2**, a benchmark for evaluating large language model agents in realistic, asynchronous environments. Unlike prior static or synchronous evaluations, Gaia2 introduces scenarios where environments evolve independently of agent actions, requiring agents to operate under temporal constraints, adapt to noisy and dynamic events, resolve ambiguity, and collaborate with other agents. Each scenario is paired with a write-action verifier, enabling fine-grained, action-level evaluation and making Gaia2 directly usable for reinforcement learning from verifiable rewards. Our evaluation of state-of-the-art proprietary and open-source models shows that no model dominates across capabilities: GPT-5 (high) reaches the strongest overall score of 42% pass@1 but fails on time-sensitive tasks, Claude-4 Sonnet trades accuracy and speed for cost, Kimi-K2 leads among open-source models with 21% pass@1. These results highlight fundamental trade-offs between reasoning, efficiency, robustness, and expose challenges in closing the “sim2real” gap. Gaia2 is built on a consumer environment with the open-source **Agents Research Environments** platform and designed to be easy to extend. By releasing Gaia2 alongside the foundational ARE framework, we aim to provide the community with a flexible infrastructure for developing, benchmarking, and training the next generation of practical agent systems.

Disclaimer

This report is **AI-GENERATED** using Large Language Models and WisPaper (a scholar search engine). It analyzes academic papers' tasks and contributions against retrieved prior work. While this system identifies **POTENTIAL** overlaps and novel directions, **ITS COVERAGE IS NOT EXHAUSTIVE AND JUDGMENTS ARE APPROXIMATE**. These results are intended to assist human reviewers and **SHOULD NOT** be relied upon as a definitive verdict on novelty.

Note that some papers exist in multiple, slightly different versions (e.g., with different titles or URLs). The system may retrieve several versions of the same underlying work. The current automated pipeline does not reliably align or distinguish these cases, so human reviewers will need to disambiguate them manually.

If you have any questions, please contact: mingzhang23@m.fudan.edu.cn

Core Task Landscape

This paper addresses: **evaluating language model agents in asynchronous dynamic environments**

A total of **50 papers** were analyzed and organized into a taxonomy with **24 categories**.

Taxonomy Overview

The research landscape has been organized into the following main categories:

- **Asynchronous and Parallel Agent Architectures**
- **Dynamic Environment Adaptation and Real-Time Decision-Making**
- **Benchmark Design and Evaluation Methodologies**
- **Sequential Planning and Reasoning Enhancement**
- **Domain-Specific Agent Applications**
- **General-Purpose Agent Frameworks and Tooling**
- **Multimodal and Context-Aware Agent Systems**
- **Specialized Technical Contributions**

Complete Taxonomy Tree

- evaluating language model agents in asynchronous dynamic environments Survey Taxonomy
- Asynchronous and Parallel Agent Architectures
 - Asynchronous Multi-Agent Coordination Frameworks (6 papers)
 - [2] Autogen: Enabling next-gen LLM applications via multi-agent conversations (Q Wu, 2024) [View paper](#)
 - [10] DynTaskMAS: A Dynamic Task Graph-driven Framework for Asynchronous and Parallel LLM-based Multi-Agent Systems (Yu Junwei, 2025) [View paper](#)
 - [12] AutoHMA-LLM: Efficient task coordination and execution in heterogeneous multi-agent systems using hybrid large language models (Tingting Yang, 2025) [View paper](#)
 - [27] MegaAgent: A large-scale autonomous LLM-based multi-agent system without predefined SOPs (Qian Wang, 2025) [View paper](#)
 - [32] Gradientsys: A Multi-Agent LLM Scheduler with ReAct Orchestration (Song Xinyuan, 2025) [View paper](#)
 - [45] Optimizing Sequential Multi-Step Tasks with Parallel LLM Agents (Zhang En-Hao, 2025) [View paper](#)
 - Single-Agent Asynchronous Systems (2 papers)
 - [24] AsyncVoice Agent: Real-Time Explanation for LLM Planning and Reasoning (Lin, 2025) [View paper](#)
 - [44] Asynchronous Tool Usage for Real-Time Agents (Kodali, 2024) [View paper](#)
 - Asynchronous RL Training and Exploration (3 papers)
 - [14] Trajectory balance with asynchrony: Decoupling exploration and learning for fast, scalable llm post-training (Bartoldson, 2025) [View paper](#)
 - [16] Beyond Ten Turns: Unlocking Long-Horizon Agentic Search with Large-Scale Asynchronous RL (Gao, 2025) [View paper](#)
 - [25] DistrL: An asynchronous distributed reinforcement learning framework for on-device control agents (Wang Taiyi, 2024) [View paper](#)
- Dynamic Environment Adaptation and Real-Time Decision-Making
 - Time-Sensitive and Rapidly Changing Environments (2 papers)
 - [7] LLM-Enhanced Rapid-Reflex Async-Reflect Embodied Agent for Real-Time Decision-Making in Dynamically Changing Environments (Mao, 2025) [View paper](#)
 - [9] Asynchronous large language model enhanced planner for autonomous driving (Yuan Chen, 2024) [View paper](#)
 - Strategic Adaptation in Multi-Agent Settings (2 papers)

- [11] Llmarena: Assessing capabilities of large language models in dynamic multi-agent environments (Junzhe Chen, 2024) [View paper](#)
- [22] Dynamic Strategy Adaptation in Multi-Agent Environments with Large Language Models (Saad, 2025) [View paper](#)
- Lifelong and Continual Learning Agents (2 papers)
- [17] "My agent understands me better": Integrating Dynamic Human-like Memory Recall and Consolidation in LLM-Based Agents (Miyashita, 2024) [View paper](#)
- [23] Lifelong Learning of Large Language Model based Agents: A Roadmap (Zheng, 2025) [View paper](#)
- Benchmark Design and Evaluation Methodologies
 - Asynchronous and Dynamic Environment Benchmarks ★ (2 papers)
 - [0] Gaia2: Benchmarking LLM Agents on Dynamic and Asynchronous Environments (Anon et al., 2026) [View paper](#)
 - [41] Inaugural MOASEI Competition at AAMAS'2025: A Technical Report (Eck, 2025) [View paper](#)
 - Domain-Specific Agent Evaluation (3 papers)
 - [1] Evaluating large language models as agents in the clinic (Nikita Mehandru, 2024) [View paper](#)
 - [5] TP-RAG: Benchmarking Retrieval-Augmented Large Language Model Agents for Spatiotemporal-Aware Travel Planning (Hang Ni, 2025) [View paper](#)
 - [20] InnovatorBench: Evaluating Agents' Ability to Conduct Innovative LLM Research (Wu Yunze, 2025) [View paper](#)
 - Task Decomposition and Tool Integration Evaluation (2 papers)
 - [15] Data Interpreter: An LLM Agent For Data Science (Sirui Hong, 2024) [View paper](#)
 - [39] Advancing Agentic Systems: Dynamic Task Decomposition, Tool Integration and Evaluation using Novel Metrics and Dataset (Ahmad, 2024) [View paper](#)
- Sequential Planning and Reasoning Enhancement
 - Search and Exploration Strategies (1 papers)
 - [19] WebRollback: Enhancing Web Agents with Explicit Rollback Mechanisms (Zhang, 2025) [View paper](#)
 - Policy Learning and Transfer (3 papers)
 - [8] A Decision-Language Model (DLM) for Dynamic Restless Multi-Armed Bandit Tasks in Public Health (Nikhil Behari, 2024) [View paper](#)
 - [29] Large Language Model as a Policy Teacher for Training Reinforcement Learning Agents (Zhou Zi-hao, 2023) [View paper](#)
 - [43] Online Intrinsic Rewards for Decision Making Agents from Large Language Model Feedback (Zheng, 2024) [View paper](#)
 - Strategic and Interactive Reasoning (2 papers)
 - [28] STRIDE: A Tool-Assisted LLM Agent Framework for Strategic and Interactive Decision-Making (Li Chuan-hao, 2024) [View paper](#)
 - [47] Plug-and-Play Policy Planner for Large Language Model Powered Dialogue Agents (Deng Yang, 2023) [View paper](#)
- Domain-Specific Agent Applications
 - Manufacturing and Industrial Control (2 papers)
 - [13] Large Language Model-Enabled Multi-Agent Manufacturing Systems (Jong-Han Lim, 2024) [View paper](#)
 - [30] A Large Language Model-Enabled Control Architecture for Dynamic Resource Capability Exploration in Multi-Agent Manufacturing Systems (Lim, 2025) [View paper](#)
 - Scientific Discovery and Engineering Design (2 papers)
 - [4] ProtAgents: protein discovery via large language model multi-agent collaborations combining physics and machine learning (Alireza Ghafarollahi, 2024) [View paper](#)
 - [34] iDesignGPT: large language model agentic workflows boost engineering design (Zhinan Zhang, 2025) [View paper](#)
 - Conversational and Interactive Agents (3 papers)
 - [37] SAUCE: Synchronous and Asynchronous User-Customizable Environment for Multi-Agent LLM Interaction (Berger, 2024) [View paper](#)
 - [38] ReSpAct: Harmonizing Reasoning, Speaking, and Acting Towards Building Large Language Model-Based Conversational AI Agents (Dongre, 2024) [View paper](#)
 - [50] Time to Talk: LLM Agents for Asynchronous Group Communication in Mafia Games (Berger, 2025) [View paper](#)
- General-Purpose Agent Frameworks and Tooling
 - Developer-Centric Agent Platforms (3 papers)
 - [18] Chainbuddy: An ai-assisted agent system for generating llm pipelines (Jingyue Zhang, 2025) [View paper](#)
 - [31] AgentScope 1.0: A Developer-Centric Framework for Building Agentic Applications (Gao Da-wei, 2025) [View paper](#)
 - [36] Agentic AI for Cloud Troubleshooting: A Review of Multi Agent System for Automated Cloud Support (Kinjal A Patel, 2025) [View paper](#)
 - Prompt Programming and LLM Integration (1 papers)
 - [40] APPL: A Prompt Programming Language for Harmonious Integration of Programs and Large Language Model Prompts (Dong Honghua, 2024) [View paper](#)
 - Hierarchical and Multi-Scale Planning (1 papers)
 - [35] Hierarchical large language model agents for multi-scale planning in dynamic environments (Yujun, 2024) [View paper](#)
- Multimodal and Context-Aware Agent Systems
 - Audio-Visual and Multimodal Understanding (1 papers)
 - [3] CAT: Enhancing Multimodal Large Language Model to Answer Questions in Dynamic Audio-Visual Scenarios (Qilang Ye, 2024) [View paper](#)
 - Data-Centric and Human-Guided Systems (1 papers)
 - [26] Towards Human-Guided, Data-Centric LLM Co-Pilots (Saveliev, 2025) [View paper](#)
 - Self-Adaptive Multi-Agent Systems (1 papers)
 - [49] Self-adaptive large language model (llm)-based multiagent systems (Nathália Nascimento, 2023) [View paper](#)
- Specialized Technical Contributions
 - Infrastructure and System Extensions (3 papers)
 - [6] Extending Oracle APEX for Large-Scale Multi-Form Workflows with Decoupled PL/SQL Logic and Asynchronous Processing Layers (Keshireddy, 2025) [View paper](#)
 - [21] An Efficient Dual-Agent Framework for Generating and Evaluating Synthetic Aviation Safety Reports using Large Language Models (Xiao Jing, 2025) [View paper](#)
 - [42] Large language model based multi-agent system augmented complex event processing pipeline (Zeeshan, 2024) [View paper](#)
 - Educational and Pedagogical Applications (2 papers)

- [33] Bringing Pedagogy into Focus: Evaluating Virtual Teaching Assistants' Question-Answering in Asynchronous Learning Environments (Li Siyan, 2025) [View paper](#)
- [46] Enhancement of online education system by using a multi-agent approach (Nethra Viswanathan, 2022) [View paper](#)
- Model Learning Dynamics and Mechanisms (1 papers)
- [48] How do language models learn facts? Dynamics, curricula and hallucinations (Zucchetti, 2025) [View paper](#)

Narrative

Core task: evaluating language model agents in asynchronous dynamic environments. The field has organized itself around several complementary perspectives. At the highest level, researchers distinguish between architectural innovations—such as asynchronous and parallel agent designs (e.g., Autogen[2], Async Planner[9])—and methodological concerns around benchmark design and evaluation (e.g., MOASEI Competition[41], Gaia2[0]). A third major branch focuses on dynamic environment adaptation and real-time decision-making, where agents must respond to shifting conditions (e.g., Dynamic Strategy Adaptation[22], Rapid-Reflex Agent[7]). Meanwhile, domain-specific applications (Clinical LLM Agents[1], ProtAgents[4]) and general-purpose frameworks (AgentScope[31], APPL[40]) reflect the tension between specialized performance and broad reusability. Additional branches address sequential planning enhancements, multimodal context integration, and specialized technical contributions, collectively mapping out a landscape that balances foundational infrastructure with task-driven innovation.

Within this ecosystem, a particularly active line of work centers on creating benchmarks that capture the unpredictability and temporal complexity of real-world settings. Gaia2[0] exemplifies this direction by emphasizing asynchronous dynamics and rigorous evaluation protocols, positioning itself alongside efforts like the MOASEI Competition[41] that stress multi-agent coordination under time pressure. In contrast, works such as CAT[3] and TP-RAG[5] prioritize context-aware reasoning and retrieval mechanisms, trading off some environmental realism for deeper semantic understanding. The interplay between these themes—whether to foreground temporal fidelity or cognitive depth—remains an open question, with Gaia2[0] leaning toward the former by foregrounding asynchronous event handling and dynamic task arrival. This choice situates it closer to benchmark-centric studies that probe agent robustness in fluid scenarios, rather than purely architectural or domain-specific explorations.

Related Works in Same Category

The following **1 sibling papers** share the same taxonomy leaf node with the original paper:

1. Inaugural MOASEI Competition at AAMAS'2025: A Technical Report

Authors: Eck, Adam, Doshi, Prashant, Soh, et al. (6 authors total) | **Year/Venue:** 2025 • arXiv.org | **URL:** [View paper](#)

Abstract

We present the Methods for Open Agent Systems Evaluation Initiative (MOASEI) Competition, a multi-agent AI benchmarking event designed to evaluate decision-making under open-world conditions. Built on the free-range-zoo environment suite, MOASEI introduced dynamic, partially observable domains with agent and task openness—settings where entities may appear, disappear, or change behavior over time. The 2025 competition featured three tracks—Wildfire, Rideshare, and Cybersecurity—each highlight...

Relationship Analysis

Both papers belong to the Asynchronous and Dynamic Environment Benchmarks category, focusing on evaluating agents in environments that evolve independently of agent actions. While Gaia2 introduces a smartphone-like consumer environment with temporal constraints, noise robustness, and multi-agent collaboration using verifiable write-action evaluation, the MOASEI Competition presents a multi-track benchmarking initiative (Wildfire, Rideshare, Cybersecurity) emphasizing open-world conditions where entities dynamically appear or disappear. The key difference is that Gaia2 provides a unified benchmark with fine-grained action-level verification for RLVR training, whereas MOASEI is a competition framework evaluating diverse participant solutions across multiple distinct domains with metrics centered on utility and robustness to perturbations.

Contributions Analysis

Overall novelty summary. The paper introduces Gaia2, a benchmark for evaluating LLM agents in asynchronous, dynamic environments, alongside the ARE (Agents Research Environments) framework and an action-level verifier. Within the taxonomy, it resides in the 'Asynchronous and Dynamic Environment Benchmarks' leaf under 'Benchmark Design and Evaluation Methodologies'. This leaf contains only two papers total, including Gaia2 itself, indicating a relatively sparse research direction. The sibling paper (MOASEI Competition) shares the focus on dynamic multi-agent scenarios but emphasizes competitive coordination rather than asynchronous temporal constraints and RL-ready verification.

The taxonomy reveals that neighboring leaves address complementary concerns: 'Domain-Specific Agent Evaluation' (healthcare, travel planning) and 'Task Decomposition and Tool Integration Evaluation' focus on specialized or multi-step reasoning without emphasizing temporal dynamics. Meanwhile, the 'Time-Sensitive and Rapidly Changing Environments' leaf under 'Dynamic Environment Adaptation' explores real-time decision-making but lacks the benchmark-centric evaluation infrastructure that Gaia2 provides. The scope note for Gaia2's leaf explicitly excludes static task benchmarks and domain-specific evaluations, positioning it at the intersection of temporal realism and general-purpose assessment.

Among 26 candidates examined across three contributions, no refutable prior work was identified. For the ARE framework (10 candidates examined, 0 refutable), the Gaia2 benchmark (10 candidates, 0 refutable), and the ARE Verifier (6 candidates, 0 refutable), the analysis found no overlapping systems that combine asynchronous environment simulation, action-level verification, and RL-ready reward signals. This suggests that within the limited search scope—focused on top-K semantic matches and citation expansion—the specific combination of features appears novel, though the search does not claim exhaustive coverage of all agent benchmarking literature.

Given the sparse population of the 'Asynchronous and Dynamic Environment Benchmarks' leaf and the absence of refutable candidates among 26 examined papers, Gaia2 appears to occupy a relatively underexplored niche. However, the limited search scope means that closely related work outside the top-26 semantic matches may exist. The analysis captures the paper's positioning within a structured taxonomy and its immediate neighborhood, but does not constitute a comprehensive survey of all agent evaluation frameworks or asynchronous simulation platforms.

This paper presents **3 main contributions**, each analyzed against relevant prior work:

Contribution 1: ARE (Agents Research Environments) framework

Description: The authors introduce ARE, a research platform providing abstractions (apps, events, notifications, scenarios) for creating simulated, asynchronous environments that evolve independently of agent actions. This framework enables reproducible benchmarking and supports reinforcement learning from verifiable rewards (RLVR).

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. Event-triggered model predictive control with deep reinforcement learning for autonomous driving

URL: [View paper](#)

Brief Assessment

Event-triggered MPC[67] focuses on event-triggered control for autonomous driving using reinforcement learning to learn trigger policies, not on creating asynchronous benchmarking environments for evaluating LLM agents. The domains and objectives are fundamentally different.

2. Deep reinforcement learning for event-driven multi-agent decision processes

URL: [View paper](#)

Brief Assessment

Event-driven Multi-Agent[73] focuses on multi-agent reinforcement learning with macro-actions in event-driven environments, but does not propose a general research platform with abstractions for creating simulated environments, reproducible benchmarking infrastructure, or support for RLVR as described in ARE.

3. Robotouille: An asynchronous planning benchmark for LLM agents

URL: [View paper](#)

Brief Assessment

Robotouille[51] focuses on asynchronous planning in a cooking simulation environment with time delays, while ARE provides a general platform for creating asynchronous, event-driven benchmarks across multiple domains with apps, notifications, and scenarios.

4. Towards spike-based machine intelligence with neuromorphic computing

URL: [View paper](#)

Brief Assessment

Spike-based Intelligence[68] focuses on neuromorphic computing and spike-based machine intelligence, not on creating research platforms for asynchronous agent environments or reinforcement learning benchmarks.

5. Asynchronous training of quantum reinforcement learning

URL: [View paper](#)

Brief Assessment

Quantum RL Training[72] focuses on asynchronous training methods for quantum reinforcement learning agents using variational quantum circuits, not on creating event-driven benchmark frameworks for agent evaluation.

6. CERiL: Continuous Event-based Reinforcement Learning

URL: [View paper](#)

Brief Assessment

CERiL[74] focuses on continuous-time reinforcement learning using event camera streams for robotic control tasks, not on creating asynchronous benchmark environments with apps, events, and notifications for evaluating LLM agents.

7. Event-based communication in distributed Q-learning

URL: [View paper](#)

Brief Assessment

Event-based Q-learning[70] focuses on reducing communication in distributed reinforcement learning systems through event-triggered control techniques, not on creating asynchronous benchmark environments with abstractions for apps, events, and scenarios as ARE does.

8. Representation learning for event-based visuomotor policies

URL: [View paper](#)

Brief Assessment

Event-based Visuomotor[75] focuses on representation learning from event camera data for visuomotor policies in robotics, not on creating asynchronous event-driven benchmarks for RL agent evaluation. The paper addresses a completely different domain (vision sensors for robot control) rather than general-purpose agent evaluation frameworks.

9. Event-Triggered Reinforcement Learning Based Joint Resource Allocation for Ultra-Reliable Low-Latency V2X Communications

URL: [View paper](#)

Brief Assessment

Event-Triggered V2X[71] focuses on resource allocation for vehicular communication networks using deep reinforcement learning, not on creating general-purpose asynchronous benchmarking platforms for evaluating LLM agents across diverse scenarios.

10. A multisynaptic spiking neuron for simultaneously encoding spatiotemporal dynamics

URL: [View paper](#)

Brief Assessment

Multisynaptic Spiking Neuron[69] focuses on spiking neural network architectures for spatiotemporal encoding in neuromorphic computing, not on agent evaluation frameworks or asynchronous event-driven benchmarks for reinforcement learning.

Contribution 2: Gaia2 benchmark

Description: The authors present Gaia2, a benchmark consisting of 1,120 human-annotated scenarios in a smartphone-like environment. It evaluates agents on capabilities including temporal awareness, adaptability to dynamic events, robustness to noise, ambiguity resolution, and multi-agent collaboration, with action-level verification suitable for RLVR training.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. Robotouille: An asynchronous planning benchmark for LLM agents

URL: [View paper](#)

Brief Assessment

Robotouille[51] is a cooking-specific benchmark testing asynchronous planning with time delays and multi-agent coordination, whereas Gaia2 evaluates agents across diverse smartphone-like apps with temporal awareness, adaptability, noise robustness, and ambiguity resolution in a consumer mobile environment.

2. Asynchronous multi-agent deep reinforcement learning under partial observability

[URL: View paper](#)

Brief Assessment

Async Partial Observability[54] focuses on multi-agent reinforcement learning with temporally extended macro-actions in robotics domains, not on benchmarking LLM agents in smartphone-like environments with human-annotated scenarios and action-level verification.

3. Temporally robust multi-agent stl motion planning in continuous time

[URL: View paper](#)

Brief Assessment

Temporally Robust STL[53] focuses on multi-agent motion planning with Signal Temporal Logic specifications in continuous time, not on benchmarking LLM agents in asynchronous smartphone-like environments with temporal reasoning and multi-agent collaboration capabilities.

4. Multi-Agent Coordination

[URL: View paper](#)

Brief Assessment

Multi-Agent Coordination[52] focuses on coordination criteria and temporal reasoning for plan analysis, not on creating a comprehensive benchmark with human-annotated scenarios, action-level verification, or RLVR training capabilities as presented in Gaia2.

5. Asynchronous multi-agent multisorted systems

[URL: View paper](#)

Brief Assessment

Multisorted Systems[58] describes threshold-based agent networks with signal propagation in continuous time, focusing on autonomous network behavior and potential dynamics. This is fundamentally different from Gaia2's smartphone environment benchmark for evaluating LLM agents on temporal awareness, dynamic events, and multi-agent collaboration with action-level verification.

6. Vaiage: A Multi-Agent Solution to Personalized Travel Planning

[URL: View paper](#)

Brief Assessment

Vaiage[55] focuses on personalized travel planning through multi-agent coordination for itinerary generation, not on benchmarking LLM agents in asynchronous environments with temporal reasoning and multi-agent collaboration capabilities.

7. Finite-Time Analysis of Asynchronous Multi-Agent TD Learning

[URL: View paper](#)

Brief Assessment

Finite-Time TD Learning[60] focuses on asynchronous multi-agent temporal difference learning for policy evaluation in reinforcement learning, not on benchmarking LLM agents in dynamic environments with human-annotated scenarios and action-level verification.

8. Dealing with interdependent activities, uncertain durations, and semantic interoperability in multi-agent plans temporal coordination.

[URL: View paper](#)

Brief Assessment

Temporal Coordination[59] focuses on temporal constraint networks for multi-agent coordination in planning domains, not on benchmarking LLM agents in smartphone-like environments with action-level verification for RLVR training.

9. Asynchronous actor-critic for multi-agent reinforcement learning

[URL: View paper](#)

Brief Assessment

Async Actor-Critic[56] focuses on asynchronous multi-agent reinforcement learning with macro-actions in robotics domains, not on benchmarking LLM agents in smartphone-like environments with temporal reasoning and dynamic events as in Gaia2.

10. TraF-Align: Trajectory-aware Feature Alignment for Asynchronous Multi-agent Perception

[URL: View paper](#)

Brief Assessment

TraF-Align[57] addresses asynchronous multi-agent perception in vehicle-to-everything (V2X) cooperative perception systems, focusing on spatial and semantic feature alignment across agents with communication delays. This is a fundamentally different domain from GAIA2's smartphone-like environment for evaluating LLM agents on temporal reasoning and multi-agent collaboration tasks.

Contribution 3: ARE Verifier for action-level evaluation

Description: The authors develop a verifier that evaluates every state-changing write action against oracle annotations, checking consistency, causality, timing, and turn-level correctness. This mechanism achieves high agreement with human annotations (0.98) and provides fine-grained credit assignment for RLVR, serving as a reusable component beyond Gaia2.

This contribution was assessed against **6 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. QEDCartographer: Automating formal verification using reward-free reinforcement learning

[URL: View paper](#)

Brief Assessment

QEDCartographer[62] focuses on formal verification in theorem proving using reinforcement learning for proof synthesis, not on action-level verification for reinforcement learning agents in dynamic environments as in the original paper.

2. "good robot! now watch this!": Repurposing reinforcement learning for task-to-task transfer

[URL: View paper](#)

Brief Assessment

Task-to-Task Transfer[63] focuses on repurposing RL models for task transfer via embedding matching, not on action-level verification mechanisms for RL evaluation. The paper does not address verifier design or fine-grained credit assignment for RLVR.

3. LaViPlan : Language-Guided Visual Path Planning with RLVR

URL: [View paper](#)

Brief Assessment

LaViPlan[65] applies RLVR to autonomous driving trajectory planning, not to general agent evaluation with action-level verification. The candidate focuses on aligning VLM reasoning with low-level trajectories in driving scenarios, whereas the original paper develops a reusable verifier for state-changing write actions across diverse agent tasks.

4. Robust Deep Reinforcement Learning Using Formal Verification

URL: [View paper](#)

Brief Assessment

Formal Verification RL[64] focuses on using formal verification methods to provide corrective signals for safe RL policy learning, not on action-level verification for credit assignment in agent benchmarking or RLVR training pipelines.

5. Leveraging Reinforcement Learning for an Efficient Windows Registry Analysis during Cyber Incident Response

URL: [View paper](#)

Brief Assessment

Windows Registry Analysis[66] focuses on cyber incident response and Windows registry analysis using reinforcement learning. The candidate's validation approach is domain-specific to cybersecurity forensics, not a general-purpose action-level verifier for LLM agents in dynamic environments.

6. SR: Teaching LLMs to Self-verify and Self-correct via Reinforcement Learning

URL: [View paper](#)

Brief Assessment

Self-verify Self-correct[61] focuses on teaching LLMs to self-verify and self-correct their own reasoning through reinforcement learning in mathematical domains, not on developing verifiers for evaluating agent actions against oracle annotations in dynamic environments.

Appendix: Text Similarity Detection

No high-similarity text segments were detected across any compared papers.

References

- [0] Gaia2: Benchmarking LLM Agents on Dynamic and Asynchronous Environments [View paper](#)
- [1] Evaluating large language models as agents in the clinic [View paper](#)
- [2] Autogen: Enabling next-gen LLM applications via multi-agent conversations [View paper](#)
- [3] CAT: Enhancing Multimodal Large Language Model to Answer Questions in Dynamic Audio-Visual Scenarios [View paper](#)
- [4] ProtAgents: protein discovery via large language model multi-agent collaborations combining physics and machine learning [View paper](#)
- [5] TP-RAG: Benchmarking Retrieval-Augmented Large Language Model Agents for Spatiotemporal-Aware Travel Planning [View paper](#)
- [6] Extending Oracle APEX for Large-Scale Multi-Form Workflows with Decoupled PL/SQL Logic and Asynchronous Processing Layers [View paper](#)
- [7] LLM-Enhanced Rapid-Reflex Async-Reflect Embodied Agent for Real-Time Decision-Making in Dynamically Changing Environments [View paper](#)
- [8] A Decision-Language Model (DLM) for Dynamic Restless Multi-Armed Bandit Tasks in Public Health [View paper](#)
- [9] Asynchronous large language model enhanced planner for autonomous driving [View paper](#)
- [10] DynTaskMAS: A Dynamic Task Graph-driven Framework for Asynchronous and Parallel LLM-based Multi-Agent Systems [View paper](#)
- [11] Llmarena: Assessing capabilities of large language models in dynamic multi-agent environments [View paper](#)
- [12] AutoHMA-LLM: Efficient task coordination and execution in heterogeneous multi-agent systems using hybrid large language models [View paper](#)
- [13] Large Language Model-Enabled Multi-Agent Manufacturing Systems [View paper](#)
- [14] Trajectory balance with asynchrony: Decoupling exploration and learning for fast, scalable llm post-training [View paper](#)
- [15] Data Interpreter: An LLM Agent For Data Science [View paper](#)
- [16] Beyond Ten Turns: Unlocking Long-Horizon Agentic Search with Large-Scale Asynchronous RL [View paper](#)
- [17] "My agent understands me better": Integrating Dynamic Human-like Memory Recall and Consolidation in LLM-Based Agents [View paper](#)
- [18] Chainbuddy: An ai-assisted agent system for generating llm pipelines [View paper](#)
- [19] WebRollback: Enhancing Web Agents with Explicit Rollback Mechanisms [View paper](#)
- [20] InnovatorBench: Evaluating Agents' Ability to Conduct Innovative LLM Research [View paper](#)
- [21] An Efficient Dual-Agent Framework for Generating and Evaluating Synthetic Aviation Safety Reports using Large Language Models [View paper](#)
- [22] Dynamic Strategy Adaptation in Multi-Agent Environments with Large Language Models [View paper](#)
- [23] Lifelong Learning of Large Language Model based Agents: A Roadmap [View paper](#)
- [24] AsyncVoice Agent: Real-Time Explanation for LLM Planning and Reasoning [View paper](#)
- [25] DistrL: An asynchronous distributed reinforcement learning framework for on-device control agents [View paper](#)
- [26] Towards Human-Guided, Data-Centric LLM Co-Pilots [View paper](#)
- [27] MegaAgent: A large-scale autonomous LLM-based multi-agent system without predefined SOPs [View paper](#)
- [28] STRIDE: A Tool-Assisted LLM Agent Framework for Strategic and Interactive Decision-Making [View paper](#)
- [29] Large Language Model as a Policy Teacher for Training Reinforcement Learning Agents [View paper](#)
- [30] A Large Language Model-Enabled Control Architecture for Dynamic Resource Capability Exploration in Multi-Agent Manufacturing Systems [View paper](#)
- [31] AgentScope 1.0: A Developer-Centric Framework for Building Agentic Applications [View paper](#)
- [32] Gradientsys: A Multi-Agent LLM Scheduler with ReAct Orchestration [View paper](#)
- [33] Bringing Pedagogy into Focus: Evaluating Virtual Teaching Assistants' Question-Answering in Asynchronous Learning Environments [View paper](#)

- [34] iDesignGPT: large language model agentic workflows boost engineering design [View paper](#)
- [35] Hierarchical large language model agents for multi-scale planning in dynamic environments [View paper](#)
- [36] Agentic AI for Cloud Troubleshooting: A Review of Multi Agent System for Automated Cloud Support [View paper](#)
- [37] SAUCE: Synchronous and Asynchronous User-Customizable Environment for Multi-Agent LLM Interaction [View paper](#)
- [38] ReSpAct: Harmonizing Reasoning, Speaking, and Acting Towards Building Large Language Model-Based Conversational AI Agents [View paper](#)
- [39] Advancing Agentic Systems: Dynamic Task Decomposition, Tool Integration and Evaluation using Novel Metrics and Dataset [View paper](#)
- [40] APPL: A Prompt Programming Language for Harmonious Integration of Programs and Large Language Model Prompts [View paper](#)
- [41] Inaugural MOASEI Competition at AAMAS'2025: A Technical Report [View paper](#)
- [42] Large language model based multi-agent system augmented complex event processing pipeline [View paper](#)
- [43] Online Intrinsic Rewards for Decision Making Agents from Large Language Model Feedback [View paper](#)
- [44] Asynchronous Tool Usage for Real-Time Agents [View paper](#)
- [45] Optimizing Sequential Multi-Step Tasks with Parallel LLM Agents [View paper](#)
- [46] Enhancement of online education system by using a multi-agent approach [View paper](#)
- [47] Plug-and-Play Policy Planner for Large Language Model Powered Dialogue Agents [View paper](#)
- [48] How do language models learn facts? Dynamics, curricula and hallucinations [View paper](#)
- [49] Self-adaptive large language model (llm)-based multiagent systems [View paper](#)
- [50] Time to Talk: LLM Agents for Asynchronous Group Communication in Mafia Games [View paper](#)
- [51] Robotouille: An asynchronous planning benchmark for LLM agents [View paper](#)
- [52] Multi-Agent Coordination [View paper](#)
- [53] Temporally robust multi-agent stl motion planning in continuous time [View paper](#)
- [54] Asynchronous multi-agent deep reinforcement learning under partial observability [View paper](#)
- [55] Vaiage: A Multi-Agent Solution to Personalized Travel Planning [View paper](#)
- [56] Asynchronous actor-critic for multi-agent reinforcement learning [View paper](#)
- [57] TraF-Align: Trajectory-aware Feature Alignment for Asynchronous Multi-agent Perception [View paper](#)
- [58] Asynchronous multi-agent multisorted systems [View paper](#)
- [59] Dealing with interdependent activities, uncertain durations, and semantic interoperability in multi-agent plans temporal coordination. [View paper](#)
- [60] Finite-Time Analysis of Asynchronous Multi-Agent TD Learning [View paper](#)
- [61] SR: Teaching LLMs to Self-verify and Self-correct via Reinforcement Learning [View paper](#)
- [62] QEDCartographer: Automating formal verification using reward-free reinforcement learning [View paper](#)
- [63] "good robot! now watch this!": Repurposing reinforcement learning for task-to-task transfer [View paper](#)
- [64] Robust Deep Reinforcement Learning Using Formal Verification [View paper](#)
- [65] LaViPlan : Language-Guided Visual Path Planning with RLVR [View paper](#)
- [66] Leveraging Reinforcement Learning for an Efficient Windows Registry Analysis during Cyber Incident Response [View paper](#)
- [67] Event-triggered model predictive control with deep reinforcement learning for autonomous driving [View paper](#)
- [68] Towards spike-based machine intelligence with neuromorphic computing [View paper](#)
- [69] A multisynaptic spiking neuron for simultaneously encoding spatiotemporal dynamics [View paper](#)
- [70] Event-based communication in distributed Q-learning [View paper](#)
- [71] Event-Triggered Reinforcement Learning Based Joint Resource Allocation for Ultra-Reliable Low-Latency V2X Communications [View paper](#)
- [72] Asynchronous training of quantum reinforcement learning [View paper](#)
- [73] Deep reinforcement learning for event-driven multi-agent decision processes [View paper](#)
- [74] CERiL: Continuous Event-based Reinforcement Learning [View paper](#)
- [75] Representation learning for event-based visuomotor policies [View paper](#)