Magentic-One & Magentic-UI

From Generalist Multi-Agent Systems to Human-in-the-Loop Agents

Yikun Han

October 2, 2025

ScienceNLP Lab Meeting

Magentic-One

Motivation: Agentic Systems Need Both Reasoning & Action

- Modern LLMs can reason; real-world tasks require tools, browsing, files, code.
- Challenges: multi-step planning, recovery from errors, long-horizon dependencies.
- Goal: A generalist system that works across domains without per-task tuning.

Background: AutoGen

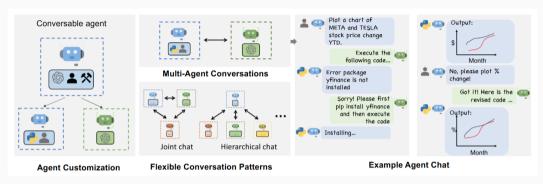


Figure 1: AutoGen enables diverse LLM-based applications using multi-agent conversations. (Left) AutoGen agents are conversable, customizable, and can be based on LLMs, tools, humans, or even a combination of them. (Top-middle) Agents can converse to solve tasks. (Right) They can form a chat, potentially with humans in the loop. (Bottom-middle) The framework supports flexible conversation patterns.

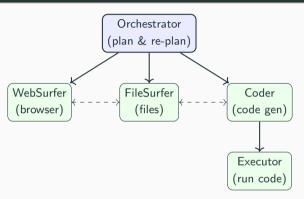
Summary

A generalist multi-agent system: an **Orchestrator** plans/re-plans and coordinates specialized agents (*WebSurfer*, *FileSurfer*, *Coder*, *Executor/Terminal*) to solve complex, open-ended tasks involving the web, files, and code.

- Modular: add/remove agents without prompt retuning.
- Competitive: GAIA, AssistantBench, WebArena.
- **Practicality**: released with *AutoGenBench* for rigorous evaluation.

System Components (High Level)

- Orchestrator: global planner, state tracker, error recovery.
- WebSurfer: navigate/search/click/extract/summarize the web.
- FileSurfer: read files, browse folders, synthesize from docs.
- **Coder**: write/analyze code; generate artifacts.
- Executor/Terminal: run code, manage env/deps.



System Components (Illustration)

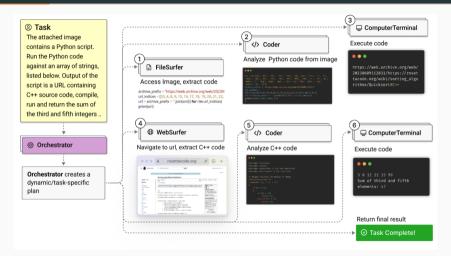


Figure 2: An illustration of the Magentic-One mutli-agent team completing a complex task from the GAIA benchmark. Magentic-One's Orchestrator agent creates a plan, delegates tasks to other agents, and tracks progress towards the goal, dynamically revising the plan as needed.

Orchestrator Agent

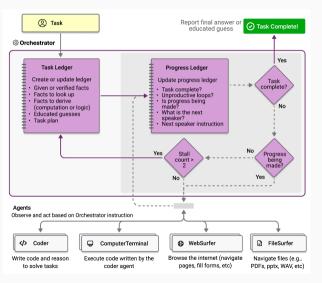


Figure 3: The Orchestrator agent that implements two loops: an outer loop and an inner loop.

Design Principles

- Generalist by default: no per-task wiring for common domains.
- Modular: plug-in agents; minimal coupling.
- **Evaluable**: AutoGenBench with repetition/isolation controls.
- Safety-aware: least privilege, logging, containment.

Orchestration & Collaboration Patterns

- $\bullet \ \ \textbf{Plan} \ \rightarrow \ \textbf{Act} \ \rightarrow \ \textbf{Observe} \ \rightarrow \ \textbf{Re-plan} \ \mathsf{loop}.$
- $\bullet \ \ \textbf{Hand-offs} \hbox{: outputs from WebSurfer} \to Coder; \ artifacts \to FileSurfer.$
- Retry & Recovery: detect stalls, adjust subgoals, switch agents.
- **State**: scratchpad, partial results, tool provenance.

Benchmarks & Result

Dataset	Category	Magentic-One (GPT-4o)	Magentic-One (GPT-4o, o1)	Best Baseline [75] [71]
GAIA [29]	Level 1 Level 2 Level 3	46.24 28.3 18.75	54.84 32.7 22.92	53.76 37.11 26.53
AssistantBench [71]	Easy Medium Hard	69.9 35.6 16.9	73.4 47.1 14.8	81 44.6 13.3
WebArena [79]	Reddit Shopping CMS Gitlab Maps Cross Site	53.77 33.16 29.1 27.78 34.86 14.6	- - - - -	65.1 36.9 24.7 39.4 33.9

Figure 4: Performance comparison between Magentic-One (GPT-4o), Magentic-One (GPT-4o, o1) and the best baseline for each benchmark's test set.

Benchmarks & Result

Dataset	Category	Magentic-One	Magentic-One	Best Baseline [75]
going east, giving the city n	ames only? Give them	to me in alphabetical order, in	a comma-separated list.	e westernmost to the easternmost
	Level 1	46.24	54.84	53.76
GAIA [29]	Level 2	28.3	32.7	37.11
	Level 3	18.75	22.92	26.53
Which supermarkets within	Easy	$\frac{1}{69.9}$ rk in Chicago have ready-to-eat	73.4	81
AssistantBench [71]	Medium	35.6	47.1	44.6
	Hard	16.9 more downvotes than upvotes f	14.8 or the user who made the lates	13.3 st post on the Showerthoughts forum
WebArena [79]	Reddit	53.77	_	65.1
	Shopping	33.16	_	36. 9
	CMS	29.1	_	24.7
	Gitlab	27.78	_	39.4
	Maps	34.86	_	33.9
	Cross Site	14.6	_	-

Figure 5: Performance comparison between Magentic-One (GPT-4o), Magentic-One (GPT-4o, o1) and the best baseline for each benchmark's test set. 10/24

Benchmarks & Result

Method	GAIA	AssistantBench (EM)	AssistantBench (accuracy)	WebArena
omne v0.1 (GPT-4o, o1)	40.53 ± 5.6	_	_	_
Trase Agent v0.2 (GPT-4o, o1,	39.53 ± 5.5	_	_	_
Gemini)				
Multi Agent (NA)	38.87 ± 5.5	_	-	_
das agent v0.4 (GPT-4o)	38.21 ± 5.5	_	_	_
Sibyl (GPT-40) [56]	34.55 ± 5.4	_	-	-
HF Agents (GPT-40)	33.33 ± 5.3	-	-	-
FRIDAY (GPT-4T) [61]	$24.25{\pm}4.8$	-	-	-
GPT-4 + plugins [29]	14.60 ± 4.0	-	-	-
$SPA \rightarrow CB (Claude)$ [71]	_	13.8 ± 5.0	26.4 ± 6.4	_
$SPA \rightarrow CB (GPT-4T) [71]$	_	9.9 ± 4.3	25.2 ± 6.3	_
Infogent (GPT-40)	_	5.5 ± 3.3	14.5 ± 5.1	_
Jace.AI (NA)	_	_	_	$57.1 {\pm} 3.4$
WebPilot (GPT-40) [75]	-	-	-	$37.2 {\pm} 3.3$
AWM (GPT-4) [57]	-	-	-	35.5 ± 3.3
SteP (GPT-4) [49]	-	-	-	33.5 ± 3.2
BrowserGym (GPT-40) [10]	_	_	-	23.5 ± 2.9
GPT-4	$6.67 \pm 2.8[29]$	$]6.1 \pm 3.5[71]$	$16.5 \pm 5.4[71]$	$14.9 \pm 2.4 [79]$
Human	92.00 ± 3.1	-	-	$78.2 {\pm} 2.8$
Magentic-One (GPT-4o)	$32.33{\pm}5.3$	11.0 ± 4.6	25.3 ± 6.3	32.8±3.2
Magentic-One (GPT-4o, o1)	38.00 ± 5.5	13.3 ± 4.9	27.7 ± 6.5	*

Figure 6: Performance of Magentic-One compared to relevant baselines on the test sets of GAIA, WebArena and AssistantBench.

Ablation Studies

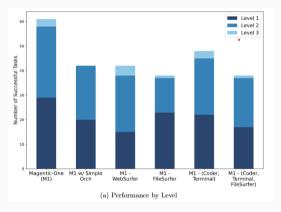


Figure 7: The performance of different ablations of Magentic-One on the GAIA validation set broken down by difficulty level.

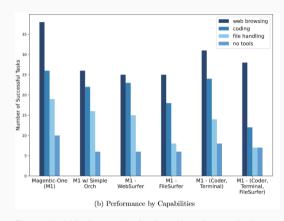


Figure 8: Ablation results broken down by required capabilities.

Error Analysis

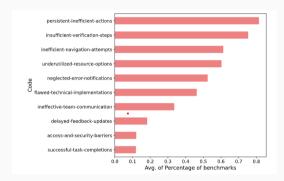


Figure 9: Distribution of error codes obtained by the automated analysis of MagenticOne's behavior as observed in the logs of the validation examples across all benchmarks studied.

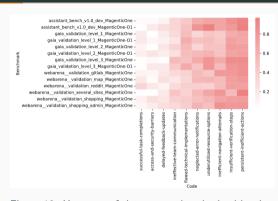


Figure 10: Heatmap of the error codes obtained by the automated analysis of Magentic-One's behavior as observed in the logs of the validation examples across all benchmarks studied.

Capabilities in Practice

- Web: search, navigate, extract tables, follow links, download.
- Files: read PDFs/CSVs/Markdown; summarize & synthesize.
- Code: prototype scripts, data processing, light automation.
- Long tasks: decompose into subgoals; checkpoint progress.

Failure Modes & Limitations

- Misplanning / Drift: wrong subgoals; fix with re-planning.
- Tool brittleness: dynamic web Uls, auth walls.
- Hallucinated affordances: attempting actions not supported.
- Safety & Security: unintended side effects; require sandboxing & oversight.

Safety & Governance Practices (Recommended)

- Least privilege: scoped credentials, read-only by default.
- Containment: containers/VMs for Executor tasks.
- Auditability: detailed logs; replay traces.
- Guardrails: domain whitelists, action confirmations on risky ops.

Extensibility & Tooling

- Implemented atop AutoGen; installable as autogen-ext[magentic-one].
- Add new agents (APIs, domain tools) with minimal prompt surgery.
- Evaluate with **AutoGenBench**; isolate side effects; repeated trials.

Magentic-UI

Why Human-in-the-Loop?

- Autonomy helps, but humans still better at judgment, risk checks.
- HITL can \(\gamma\) reliability, \(\gamma\) safety, and \(\gamma\) sample efficiency.
- ullet Magentic-UI: a web UI + multi-agent backend that operationalizes HITL.

What is Magentic-UI?

- Open-source research prototype.
- Built on a flexible multi-agent stack derived from Magentic-One.
- Supports web browsing, file ops, code gen/exec; extensible via MCP.

What is Magentic-UI?

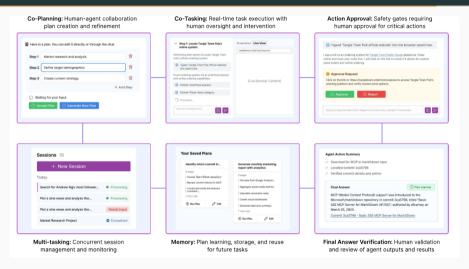


Figure 11: agentic-UI is an open-source research prototype of a human-centered agent that is meant to help researchers study open questions on human-in-the-loop approaches and oversight.

HITL Mechanisms in Magentic-UI

- **Co-planning**: user edits the plan before execution.
- Co-tasking & multi-tasking: parallel goals with human arbitration.
- Action guards: confirm/deny sensitive actions.
- Long-term memory: reuse context across sessions.

Building on Magentic-One

- Reuses the **Orchestrator** + **Specialists** pattern.
- Adds UI hooks for plan steps, tool outputs, and confirmations.
- Emphasizes **transparency**: show actions, states, artifacts.

When to Prefer HITL

- High-risk actions (payments, emails, system changes).
- Ambiguous specs or noisy environments.
- \bullet Novel tasks where autonomy struggles; HITL \Rightarrow faster iteration.

Discussion

Key Takeaways

- Magentic-One: generalist, modular, evaluated, practical.
- Magentic-UI: human-centered controls on top of the same stack.
- For research groups, combine autonomy with **HITL** to move fast safely.

References

References i



A. Fourney et al.

Magentic-One: A Generalist Multi-Agent System for Solving Complex Tasks.

arXiv:2411.04468, 2024.

https://arxiv.org/abs/2411.04468



Microsoft Research Blog.

Magentic-One: A Generalist Multi-Agent System for Solving Complex Tasks.

2024.

https://www.microsoft.com/en-us/research/articles/ magentic-one-a-generalist-multi-agent-system-for-solving-complex-tasks/

References ii

AutoGen Docs.

Magentic-One user guide & API.

https://microsoft.github.io/autogen/stable//user-guide/agentchat-user-guide/magentic-one.html

H. Mozannar et al.

Magentic-UI: Towards Human-in-the-loop Agentic Systems.

arXiv:2507.22358, 2025.

https://arxiv.org/abs/2507.22358

References iii

Microsoft Research Blog.

Magentic-UI, an experimental human-centered web agent. 2025.

https://www.microsoft.com/en-us/research/blog/magentic-ui-an-experimental-human-centered-web-agent/

Microsoft GitHub.

Magentic-UI repository.

https://github.com/microsoft/magentic-ui