# W2.2- LLM Agent Basics

2026 Spring
LLM Agents Foundation & Applications

Dr. Yanjun Qi

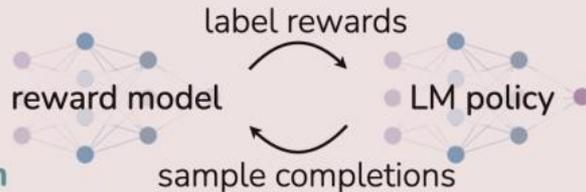20260115

# Last Class

To incorporate human preferences:

- **RL algorithms** (PPO, GRPO, REINFORCE) explicitly maximize expected reward from a reward model – we normally call this group **RLHF**
- **Direct alignment methods** (DPO, IPO, KTO) optimize preference objectives without explicit reward modeling
- Though they can be shown to implicitly optimize an equivalent objective under certain assumptions

# Accelerated development (25 Sept)

Text input → LLM → Text output

Parameters (Bn)   open access

## Major Large Language Models (LLMs)
ranked by capabilities, sized by billion parameters used for training

CLICK LEGEND ITEMS TO FILTER

● anthropic ● chinese ● google ● meta ● mistral ● openAI ● other ● xAI

search...    show only: all



MMLU

89.8 = human expert

80

▲ 70+ IDEAL ▲

60

40

20

pre-2022    2022    2023    2024    2025

David McCandless, Tom Evans, Paul Barton
**Informationisbeautiful** // Sep 2025

MMLU = benchmark for measuring LLM capabilities
* = parameters undisclosed // source: LifeArchitect // data

3

# LLM agents: enabling LLMs to interact with the environment

# LLM Agents in Diverse Environments

# Multi-agent collaboration: division of labor for complex tasks



**Emergence of social behaviors with role-play LLMs**
Generative agents, Project Sid,...

**Specialized agents for different subtasks**
Autogen, CrewAI, CAMEL, Mixture-of-Agents,...

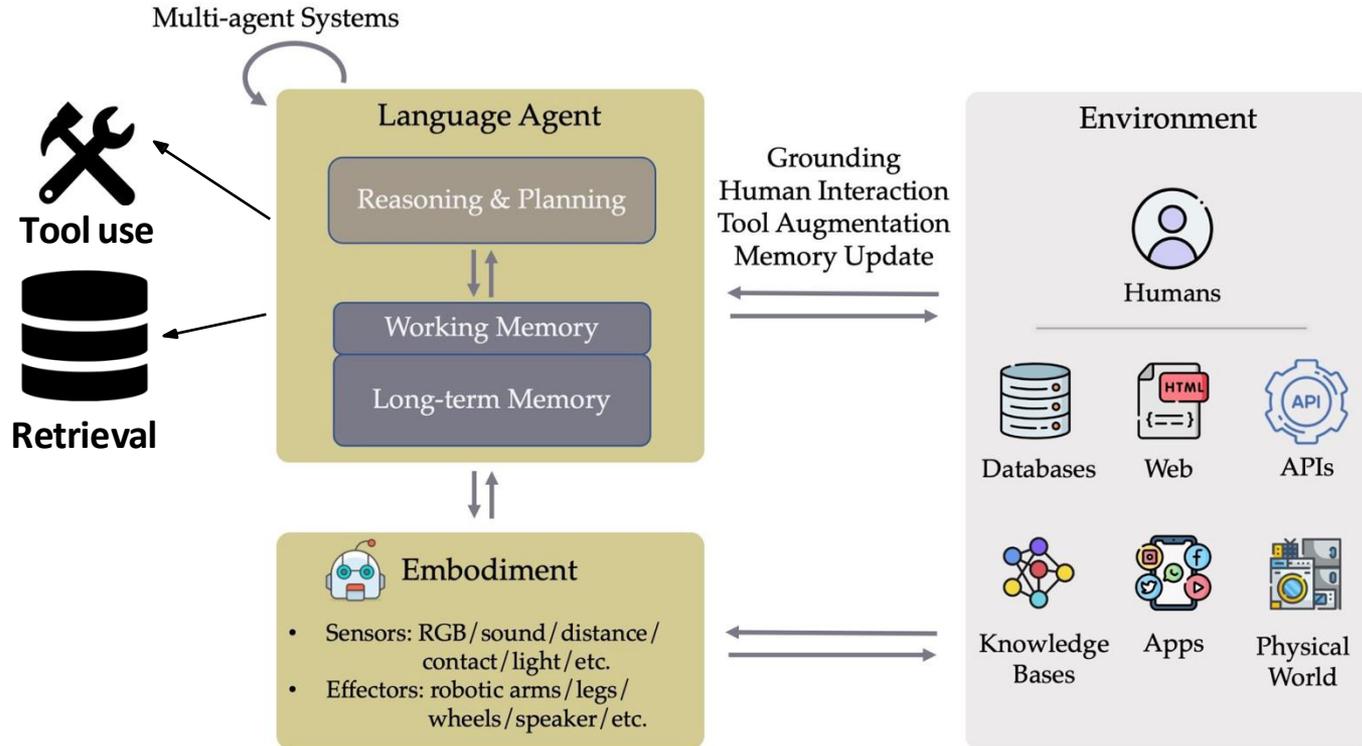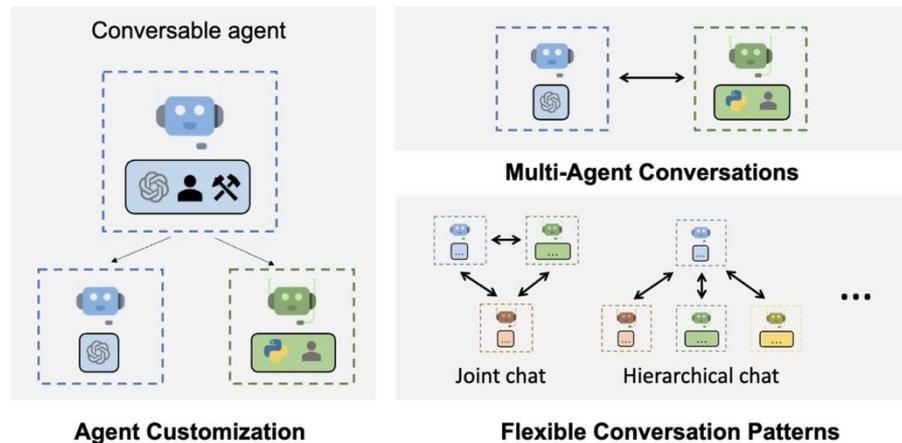Project Sid: Many-agent simulations toward AI civilization, 2024

# Why empowering LLMs with the agent framework

**Memory**

**Tool use**

**Retrieval**

**LLM**
**Reasoning & Planning**

**Agent**

- Solving real-world tasks typically involves a trial-and-error process

- Leveraging external tools and retrieving from external knowledge expand LLM's information capabilities

- Agent workflow facilitates complex tasks
  - Allocation of subtasks to specialized tools
  - Multi-agent generation inspires better responses
  - Access to specialized evidence / data / inputs
  - ......

# LLM agents transformed various applications



**Code generation**
Cursor, GitHub Copilot, Devin, Replit,…



**Workflow automation**
Microsoft Copilot, Multi-On,…



**Personal assistant**
Google Astra, OpenAI GPT-4o,…



**Robotics**
Figure AI, Tesla Optimus,…

- Healthcare
- Education
- Law
- Finance
- Cybersecurity

    …

8

# LLM agents are improving rapidly (leaderboards!^

SWE-bench

*Full* is a large benchmark made of 2000 instances (details)

Filters: Open Scaffold ▾    All Tags ▾

| Model | % Resolved | Org |
|---|---|---|
| ✅ SWE-agent 1.0 (Claude 3.7 Sonnet) | 33.83 | |
| ✅ OpenHands + CodeAct v2.1 (claude-3-5-sonnet-20241022) | 29.38 | |
| AutoCodeRover-v2.0 (Claude-3.5-Sonnet-20241022) | 24.89 | |
| ✅ SWE-agent + Claude 3.5 Sonnet | 18.13 | |
| ✅ SWE-agent + GPT 4 (1106) | 12.47 | |
| ✅ SWE-agent + GPT 4o (2024-05-13) | 11.99 | |
| ✅ SWE-agent + Claude 3 Opus | 10.51 | |
| ✅ RAG + Claude 3 Opus | 3.79 | |

Leaderboards

**BENCHMARKS**

SWE-bench
SWE-bench Verified ⧉
SWE-bench Bash Only
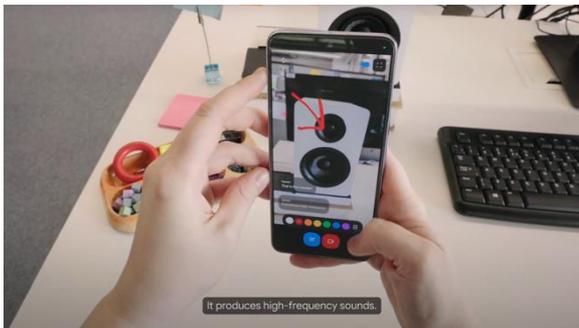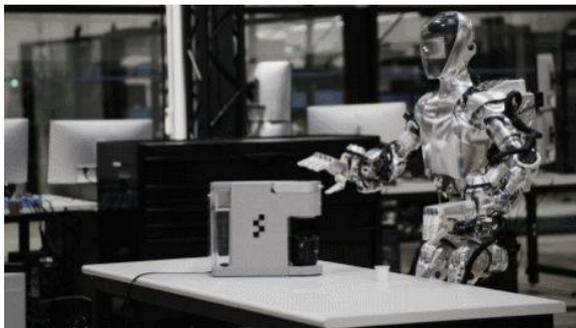SWE-bench Multilingual
SWE-bench Multimodal
SWE-bench Lite

**ABOUT**

SWE-Bench (Jimenez*, Yang*, et al.) / swebench.com / Today's screenshot on Leaderboard!

GAIA (Mialon et al.)
huggingface.co/gaia-benchmark

Results: Test

| Agent name | Model family |
|---|---|
| JoinAI V2.2 | GPT 5, Gemini 3 Pro, DeepSeek 3.1, Qwen 3 |
| Nemotron-ToolOr | Nemotron-ToolOrchestrator-8B, GPT-5, Claude Opus 4.1 |
| Nemotron-ToolOr | Nemotron-ToolOrchestrator-8B, GPT-5, Claude Opus 4.1 |
| SU Zero - Shuqi | Self Consistency 35 |
| JoinAI V2.1 | GPT, Gemini, DeepSeek, Qwen |
| ShawnAgent v3.1 | GPT5.2, Claude Sonnet 4.5, Gemini 3 Pro |
| HALO V1217-1 | |

F2    OAgent

| | A | B | C | D | E | F | G | H | |
|---|---|---|---|---|---|---|---|---|---|
| 1 | a | Open? | Model Size (billion) | Model | Success Rate (%) | Result Source | Work | Traj | |
| 2 | 01/2026 | ✗ | | OAgent | 71.6 | OAgent | OAgent | Link | |
| 3 | 12/2025 | ✔ | GPT-5 | ColorBrowserAgent | 71.2 | ColorBrowserAgent | ColorBrowserAgent | Link | ite-s |
| 4 | 10/2025 | ✔ | - | Claude Code + GBOX MCP | 68 | GBOX AI | GBOX AI | Link | |
| 5 | 09/2025 | ✗ | - | DeepSky Agent | 66.9 | Self-reported | DeepSky Agent | Link | |
| 6 | 10/2025 | ✗ | - | Narada AI | 64.2 | Self-reported | Narada AI | Link | |
| 7 | 02/2025 | ✔ | - | IBM CUGA | 61.7 | IBM CUGA | IBM CUGA | html+ json | |
| 8 | 01/2025 | ✗ | - | OpenAI Operator | 58.1 | OpenAI CUA | OpenAI CUA | Link | Syste |

WebArena (Zhou et al.)
webarena.dev

# Frontier performance across benchmarks

**Accuracy**

53 Results

**Benchmark**

- Mock AIME 24-25
- GPQA Diamond
- FrontierMath Tier 4
- SimpleQA Verified
- SWE-bench Verified



GPT-5.2 (high)
Gemini 3 Pro Preview
o1 (high)
Gemini 3 Pro Preview
o3 (high)
o1-mini (high)
GPT-4 (Mar 2023)
Claude 3.5 Sonnet (Jun 2024)
GPT-5.2 (Pro)
GPT-5 (high)
GPT-4 (Mar 2023)
Claude 3.5 Sonnet (Jun 2024)

**Release Date**

Apr. 2023  July 2023  Oct. 2023  Jan. 2024  Apr. 2024  July 2024  Oct. 2024  Jan. 2025  Apr. 2025  July 2025  Oct. 2025  Jan. 2026

# Agent architecture



Image Credit: Percy Liang from Stanford U.

# Topics We will cover in this course:

# Topics covered in this course. ➡ Potential Projects:

- Applications
  - Software development
  - Workflow automation
  - Multimodal applications
  - Industrial applications like Healthcare, Legal, Fin, …
- Model core capabilities
  - Reasoning
  - Planning
  - Multimodal understanding
- LLM agent frameworks
  - Workflow
  - Tool use
  - Retrieval-augmented generation
  - Multi-agent systems
- Safety and ethics

- Applications
  - Build LLM agents applications in specialized / novel domains
- Core Fundamentals
  - Enhance core agent capabilities (memory, planning, tool use, alignments, efficiency, …)
  - Enhance decentralized multi-agent systems
- Benchmarks / Build Novel Frameworks
  - Create and improve benchmarks for Evaluating LLM agents
  - Reimplement or Build novel frameworks for agent workflow
- Safety and ethics
  - Reveal safety concerns in deployment (misuse, privacy, etc.)
  - Defense safety converns

13

# Challenges for LLM agent deployment in the wild

- ## Reasoning and planning
  - LLM agents tend to make mistakes when performing complex tasks end-to-end

- ## Embodiment and learning from environment feedback
  - LLM agents are not yet efficient at recovering from mistakes for long-horizon tasks
  - Continuous learning, self-improvement
  - Multimodal understanding, grounding and world models

- ## Multi-agent learning, theory of mind

- ## Safety and privacy
  - LLMs are susceptible to adversarial attacks, can emit harmful messages and leak private data

- ## Human-agent interaction, ethics
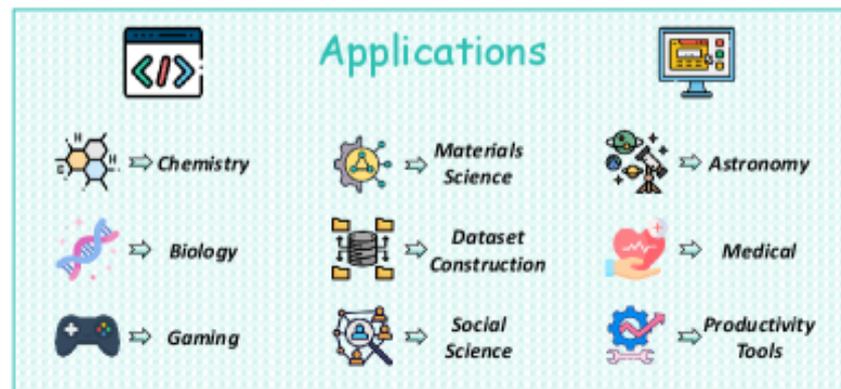  - How to effectively control the LLM agent behavior, and design the interaction mode between humans and LLM agents
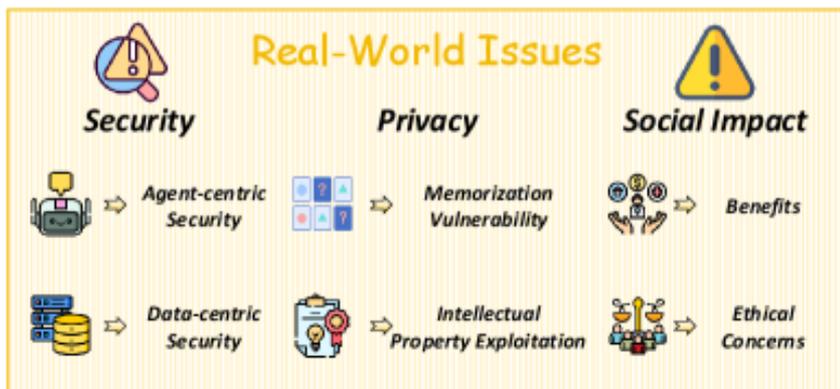
# Large Language Model Agent: A Survey on Methodology, Applications and Challenges

Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao, Junwei Yang, Yiyang Gu, Bohan Wu, Binqi Chen,
Ziyue Qiao, Qingqing Long, Rongcheng Tu, Xiao Luo, Wei Ju, Zhiping Xiao, Yifan Wang, Meng Xiao,
Chenwu Liu, Jingyang Yuan, Shichang Zhang, Yiqiao Jin, Fan Zhang, Xian Wu, Hanqing Zhao,
Dacheng Tao, *Fellow, IEEE*, Philip S. Yu, *Fellow, IEEE* and Ming Zhang

**Abstract**—The era of intelligent agents is upon us, driven by revolutionary advancements in large language models. Large Language Model (LLM) agents, with goal-driven behaviors and dynamic adaptation capabilities, potentially represent a critical pathway toward artificial general intelligence. This survey systematically deconstructs LLM agent systems through a methodology-centered taxonomy, linking architectural foundations, collaboration mechanisms, and evolutionary pathways. We unify fragmented research threads by revealing fundamental connections between agent design principles and their emergent behaviors in complex environments. Our work provides a unified architectural perspective, examining how agents are constructed, how they collaborate, and how they evolve over time, while also addressing evaluation methodologies, tool applications, practical challenges, and diverse application domains. By surveying the latest developments in this rapidly evolving field, we offer researchers a structured taxonomy for understanding LLM agents and identify promising directions for future research. The collection is available at https://github.com/luo-junyu/Awesome-Agent-Papers.

**Index Terms**—Large language model, LLM agent, AI agent, intelligent agent, multi-agent system, LLM, literature survey

27 Mar 2025

15

# Agent Methodology

### Construction
- Profile Definition
- Memory Mechanism
- Planning Capability
- Action Execution

### Collaboration
- Centralized Control
- Decentralized Collaboration
- Hybrid Architecture

### Evolution
- Self-Learning
- Multi-agent Co-Evolution
- External Resource

# Evaluation and Tools

### Benchmark and Datasets
- General Assessment
- Domain-specific Evaluation
- Collaboration Evaluation

### Tools
- LLM Use Tools
- LLM Create Tools
- Tools Develop LLM

# Real-World Issues

### Security
- Agent-centric Security
- Data-centric Security

### Privacy
- Memorization Vulnerability
- Intellectual Property Exploitation

### Social Impact
- Benefits
- Ethical Concerns

# Applications
- Chemistry
- Materials Science
- Astronomy
- Biology
- Dataset Construction
- Medical
- Gaming
- Social Science
- Productivity Tools

16

```
                    ┌─────────────────────┐        ┌──────────────────────────────────────────────────────────────┐
                    │  Human-Curated      │        │ Camel [25], AutoGen [26], MetaGPT [27], ChatDev [28], AFlow [29] │
    ┌───────────────┤  Static Profiles    ├────────┤                                                                │
    │ Profile       │                     │        └──────────────────────────────────────────────────────────────┘
    │ Definition    │                     │
    │ §2.1.1        │        ┌─────────────────────┐        ┌──────────────────────────────────────────────────────┐
    └───────────────┤  Betch-Generated    ├────────┤ Generative Agents [30], RecAgent [31], DSPy [32]       │
                    │  Dynamic Profiles   │        │                                                        │
                    └─────────────────────┘        └──────────────────────────────────────────────────────┘

                    ┌─────────────────────┐        ┌──────────────────────────────────────────────────────────────┐
                    │  Short-Term         │        │ ReAct [33], ChatDev [28], Graph of Thoughts [34], AFlow [29]   │
    ┌───────────────┤  Memory             ├────────┤                                                                │
    │ Memory        │                     │        └──────────────────────────────────────────────────────────────┘
    │ Mechanism     │        ┌─────────────────────┐        ┌──────────────────────────────────────────────────────────────┐
    │ §2.1.2        │        │  Long-Term          │        │ Voyager [35], GITM [36], ExpeL [37], Reflexion [38], TPTU [39], │
    └───────────────┤  Memory             ├────────┤ OpenAgents [40], Lego-Prover [41], MemGPT [42]                  │
                    │                     │        └──────────────────────────────────────────────────────────────┘
                    │        ┌─────────────────────┐        ┌──────────────────────────────────────────────────────────────┐
                    │        │  Knowledge          │        │ RAG [43], GraphRAG [44], Chain of Agnets [45], IRCoT [46],      │
                    └────────┤  Retrieval          ├────────┤ Llatrieval [47], KG-RAR [48], DeepRAG [49]                      │
                             │  as Memory          │        └──────────────────────────────────────────────────────────────┘
                             └─────────────────────┘
```

Language Model Agent

**Planning Capability §2.1.3**

- **Task Decomposition Strategies** → Plan-and-solve Prompting [50], Distributed Problem Solving and Planning [51], ReAct [33], Chain-of-discussion [52], Tree-planner [53], ReAcTree [54], ToT [55], ReST-MCTS* [56], LLM-MARS [57], LLM as BT-planner [58], ConceptAgent [59]

- **Feedback-Driven Iteration** → BrainBody-LLM [60], TrainerAgent [61], RASC [62], REVECA [63], AdaPlanner [64], AIFP [65]

**Action Execution §2.1.4**

- **Tool Utilization** → TRICE [66], GPT4Tools [67], EASYTOOL [68], AvaTaR [69],

- **Physical Interaction** → DriVLMe [70], ReAd [71], Collaborative Voyager [72]

```
                    ┌─────────────────┐  ┌──────────────────────────────────────────────────┐
                    │ Centralized     │  │ Coscientist [73], LLM-Blender [74], MetaGPT [27], │
                    │ Control         │──│ AutoAct [75],                                     │
                    │ §2.2.1          │  │ Meta-Prompting [76], Wjudge [77]                  │
                    └─────────────────┘  └──────────────────────────────────────────────────┘
   ┌──────────────┐ ┌─────────────────┐  ┌──────────────────────────────────────────────────┐
   │ Agent        │ │ Decentralized   │  │ GAgents [30], CAMEL [25], MedAgents [78],         │
   │ Collaboration│─│ Collaboration   │──│ ReConcile [79], MAD [80], MADR [81], MDebate [82],│
   │ §2.2         │ │ §2.2.2          │  │ AutoGen [26]                                      │
   └──────────────┘ └─────────────────┘  └──────────────────────────────────────────────────┘
                    ┌─────────────────┐  ┌──────────────────────────────────────────────────┐
                    │ Hybrid          │  │                                                  │
                    │ Architecture    │──│ KnowAgent [83], WKM [84], Textgrad [85]          │
                    │ §2.2.3          │  │                                                  │
                    └─────────────────┘  └──────────────────────────────────────────────────┘

                    ┌─────────────────┐  ┌──────────────────────────────────────────────────┐
                    │ Autonomous      │  │ SE [86], Evolutionary Optimization [87],         │
                    │ Optimization and│  │ DiverseEvol [88], SELF-REFINE [89], STaR [90],   │
                    │ Self-Learning   │──│ V-STaR [91], Self-Verification [92],             │
                    │ §2.3.1          │  │ Self-Rewarding [93], RLCD [94], RLC [95]         │
                    └─────────────────┘  └──────────────────────────────────────────────────┘
   ┌──────────────┐ ┌─────────────────┐  ┌──────────────────────────────────────────────────┐
   │ Agent        │ │ Multi-Agent     │  │ ProAgent [96], CORY [97], CAMEL [25],            │
   │ Evolution    │─│ Co-Evolution    │──│ Red-Team LLMs [98],                             │
   │ §2.3         │ │ §2.3.2          │  │ Multi-Agent Debate [82], MAD [99]               │
   └──────────────┘ └─────────────────┘  └──────────────────────────────────────────────────┘
                    ┌─────────────────┐  ┌──────────────────────────────────────────────────┐
                    │ Evolution via   │  │ KnowAgent [83], WKM [84], CRITIC [100],          │
                    │ External        │──│ STE [101], SelfEvolve [102]                      │
                    │ Resources §2.3.3│  │                                                  │
                    └─────────────────┘  └──────────────────────────────────────────────────┘
```

# TABLE 2: A summary of agent evolution methods.

| Category | Method | Key Contribution |
|---|---|---|
| **Self-Supervised Learning** | SE [86] | Adaptive token masking for pretraining |
| | Evolutionary Optimization [87] | Efficient model merging and adaptation |
| | DiverseEvol [88] | Improved instruction tuning via diverse data |
| **Self-Reflection & Self-Correction** | SELF-REFINE [89] | Iterative self-feedback for refinement |
| | STaR [90] | Bootstrapping reasoning with few rationales |
| | V-STaR [91] | Training a verifier using DPO |
| | Self-Verification [92] | Backward verification for correction |
| **Self-Rewarding & RL** | Self-Rewarding [93] | LLM-as-a-Judge for self-rewarding |
| | RLCD [94] | Contrastive distillation for alignment |
| | RLC [95] | Evaluation-generation gap for optimization |
| **Cooperative Co-Evolution** | ProAgent [96] | Intent inference for teamwork |
| | CORY [97] | Multi-agent RL fine-tuning |
| | CAMEL [25] | Role-playing framework for cooperation |
| **Competitive Co-Evolution** | Red-Team LLMs [98] | Adversarial robustness training |
| | Multi-Agent Debate [82] | Iterative critique for refinement |
| | MAD [99] | Debate-driven divergent thinking |
| **Knowledge-Enhanced Evolution** | KnowAgent [83] | Action knowledge for planning |
| | WKM [84] | Synthesizing prior and dynamic knowledge |
| **Feedback-Driven Evolution** | CRITIC [100] | Tool-assisted self-correction |
| | STE [101] | Simulated trial-and-error for tool learning |
| | SelfEvolve [102] | Automated debugging and refinement |

Real-world Issues

**Security** — Agent-centric Security: Adversarial Attacks, Jailbreak Attacks, Backdoor Attacks, Model Collaboration. Data-centric Security: External Data Attacks, Interaction Attacks.

**Privacy** — Memorization Vulnerability: Data Extraction Attacks, Member Inference Attacks, Attribute Inference Attacks. Intellectual Property Exploitation: Model Stealing Attacks, Prompt Stealing Attacks.

**Social Impact** — Benefits: Automation Enhancement, Job Creation and Workforce Transformation, Enhance Information Distribution. Ethical Concerns: Bias and Discrimination, Accountability, ...

**TABLE 3: Summary of agent-centric attacks and defense in LLM agents.**

| Reference | Description |
|---|---|
| **Adversarial Attacks and Defense** | |
| Mo et al. [177] | **Attack:** Adversarial attack benchmark |
| AgentDojo [178] | **Attack:** Adversarial attack framework |
| ARE [179] | **Attack:** Adversarial attack evaluation for multimodal agents |
| GIGA [181] | **Attack:** Generalizable infectious gradient attacks |
| CheatAgent [180] | **Attack:** Adversarial attack agent for recommender systems |
| LLAMOS [182] | **Defense:** Purifying adversarial attack input |
| Chern et al. [183] | **Defense:** Defense via multi-agent debate |
| **Jailbreaking Attacks and Defense** | |
| RLTA [184] | **Attack:** Produce jailbreaking prompts via reinforcement learning |
| Atlas [185] | **Attack:** Jailbreaks text-to-image models with safety filters |
| RLbreaker [186] | **Attack:** Model jailbreaking as a search problem |
| PathSeeker [187] | **Attack:** Use multi-agent reinforcement learning to jailbreak |
| AutoDefense [188] | **Defense:** Multi-agent defense to filter harmful responses |
| Guardians [189] | **Defense:** Detect rogue agents to counter jailbreaking attacks. |
| ShieldLearner [190] | **Defense:** Learn attack jailbreaking patterns. |
| **Backdoor Attacks and Defense** | |
| DemonAgent [191] | **Attack:** Encrypted muti-backdoor implantation attack |
| Yang et al. [192] | **Attack:** Backdoor attacks evaluations on LLM-based agents |
| BadAgent [193] | **Attack:** Inputs or environment cues as backdoors |
| BadJudge [194] | **Attack:** Backdoor to the LLM-as-a-judge agent system |
| DarkMind [195] | **Attack:** latent backdoor attack to customized LLM agents |
| **Agent Collaboration Attacks and Defense** | |
| CORBA [196] | **Attack:** Multi-agent attack via multi-agent |

TABLE 4: Summary of data-centric attack and defense in LLM agents.

| Reference | Description |
|---|---|
| **External Data Attacks and Security** | |
| Li et al. [204] | **Attack:** Malicious prefix injection |
| Psysafe [201] | **Attack:** A dark psychological injection benchmark |
| Tian et al. [210] | **Attack:** Guide agents into specific role-playing states |
| InjectAgent [205] | **Attack:** A prompting injection benchmark |
| Agentdojo [203] | **Attack:** A user injection benchmark |
| AgentPoison [216] | **Attack:** Poisoning samples in knowledge databases |
| Nakash et al. [215] | **Attack:** Indirect prompt injection through FITD attack |
| WIPI [214] | **Attack:** control agents through a public web page |
| ASB [176] | **Attack:** A multi-type attack benchmark |
| AgentHarm [223] | **Attack:** A multi-type attack benchmark |
| Mantis [206] | **Defense:** Hacking back to attackers |
| Chern et al. [183] | **Defense:** Employ multi-agent debate to verify external knowledge |
| RTBAS [208] | **Defense:** Check every step of agent information flow |
| TaskShield [209] | **Defense:** Check every step of agent process |
| Zhang et al. [201] | **Defense:** Doctor and police agents guard the healthy psychology |
| **Interaction Attacks and Security** | |
| Wang et al. [217] | **Attack:** Private memory extraction attack |
| CORBA [196] | **Attack:** Disrupt the communications among agents |
| AgentSmith [220] | **Attack:** Poison one agent to infectious other agents |
| Lee et al. [221] | **Attack:** Conduct injections to self-replicate among agents |
| He et al. [197] | **Attack:** Inject semantic disruptions to agent communications |
| BlockAgents [222] | **Defense:** Incorporate blockchain and PoT against byzantine attacks |
| Abdelnabi et al. [207] | **Defense:** A multi-layer agent firewall |

## TABLE 7: Overview of Applications in LLM Agents.

| Method | Domain | Core Idea |
|---|---|---|
| **Scientific Discovery** | | |
| SciAgents [266] | General Sciences | Collaborative hypothesis generation |
| Curie [267] | General Sciences | Automated experimentation |
| ChemCrow [269] | Chemistry | Tool-augmented synthesis planning |
| AtomAgents [270] | Materials Science | Physics-aware alloy design |
| D. Kostunin el al [271] | Astronomy | Telescope configuration management |
| BioDiscoveryAgent [273] | Biology | Genetic perturbation design |
| GeneAgent [274] | Biology | Self-verifying gene association discovery |
| RiGPS [275] | Biology | Biomarker identification |
| BioRAG [211] | Biology | Biology-focused retrieval augmentation |
| PathGen-1.6M [276] | Medical Dataset | Pathology image dataset generation |
| KALIN [277] | Biology Dataset | Scientific question corpus generation |
| GeneSUM [278] | Biology Dataset | Gene function knowledge maintenance |
| AgentHospital [281] | Medical | Virtual hospital simulation |
| ClinicalLab [282] | Medical | Multi-department diagnostics |
| AIPatient [283] | Medical | Patient simulation |
| CXR-Agent [284] | Medical | Chest X-ray interpretation |
| MedRAX [285] | Medical | Multimodal medical reasoning |
| **Gaming** | | |
| ReAct [33] | Game Playing | Reasoning and acting in text environments |
| Voyager [35] | Game Playing | Lifelong learning in Minecraft |
| ChessGPT [287] | Game Playing | Chess gameplay evaluation |
| GLAM [288] | Game Playing | Reinforcement learning in text environments |
| CALYPSO [289] | Game Generation | Narrative generation for D&D |
| GameGPT [290] | Game Generation | Automated game development |
| Sun et al. [291] | Game Generation | Interactive storytelling experience |

# A Survey on Large Language Model based Autonomous Agents

Lei Wang, Chen Ma*, Xueyang Feng*, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhi-Yuan Chen, Jiakai Tang, Xu Chen(✉), Yankai Lin(✉), Wayne Xin Zhao, Zhewei Wei, Ji-Rong Wen

Gaoling School of Artificial Intelligence, Renmin University of China, Beijing, 100872, China

**Abstract** Autonomous agents have long been a research focus in academic and industry commu-

LLM-based autonomous agents. Based on the pre-

**Number of Papers（cumulated）**

**Time（Year-Month）**

Legend:
- General Agent
- Tool Agent
- Simulation Agent
- Embodied Agent
- Game Agent
- Web Agent
- Assistant Agent

WebGPT 2021-12

CoT 2022-1

TALM 2022-5

WebShop 2022-7

Inner Monologue 2022-7

DEPS 2023-2

Toolformer 2023-2

AutoGPT 2023-3

HuggingGPT 2023-3

Generative Agent 2023-4

AgentGPT 2023-4

Voyager 2023-5

GITM 2023-5

ToT 2023-5

MIND2WEB 2023-6

RecAgent 2023-6

ChatDev 2023-7

CO-LLM 2023-7

ToolBench 2023-7

AgentSims 2023-8

## Profile

**Profile Contents**
- Demographic Information
- Personality Information
- Social Information

**Generation Strategy**
- Handcrafting Method
- LLM-Generation Method
- Dataset Alignment Method

## Memory

**Memory Structure**
- Unified Memory
- Hybrid Memory

**Memory Formats**
- Languages
- Databases
- Embeddings
- Lists

**Memory Operation**
- Memory Reading
- Memory Writing
- Memory Reflection

## Planning

**Planning w/o Feedback**
- Single-path Reasoning
- Multi-path Reasoning
- External Planner

**Planning w/ Feedback**
- Environment Feedback
- Human Feedback
- Model Feedback

## Action

**Action Target**
- Task Completion
- Exploration
- Communication

**Action Production**
- Memory Recollection
- Plan Following

**Action Space**
- Tools
- Self-Knowledge

**Action Impact**
- Environments
- New Actions
- Internal States

**CoT , Zero-shot Cot**

Prompts

LLM

Reasoning Step-1

Reasoning Step-2

Reasoning Step-n

**ReWOO , HuggingGPT**

Prompts

LLM

Reasoning Step-1

LLM

Reasoning Step-2

LLM

Reasoning Step-n

**Single-Path Reasoning**

**CoT-SC**

Prompts

LLM

| Step-1 | Step-1 | Step-1 |
| Step-2 | Step-2 | Step-2 |
| Step-n | Step-n | Step-n |

**ToT , LMZSP , RAP**

Prompts

LLM

Step-1

Step-2 | Step-2 | Step-2

Step-3 | Step-3 | Step-3 | Step-3

**Multi-Path Reasoning**

**Parameter Learning**

Dataset

Output

Model

Input

Model Parameters

Capability

**Parameter Learning**

Model

**Prompt Engineering**

Classify the text into neutral, negative or positive. Text: I think the food was okay. Sentiment:

Output

Neutral

Prompts

Capability

**Parameter Learning**

Model

**Prompt Engineering**

Agent relevant Prompts

Output

**Mechanism Engineering**

Trial-and-Error

Crowd-sourcing

MECHANISMS

Capability

The era of machine learning

The era of large language model

The era of agent

**Table 2**    Representative applications of LLM-based autonomous agents.

| | Domain | Work |
|---|---|---|
| Social Science | Psychology | TE [101], Akata et al. [102], Ziems et al. [104], Ma et al. [103] |
| | Political Science and Economy | Argyle et al. [29], Horton [105], Ziems et al. [104] |
| | Social Simulation | Social Simulacra [80], Generative Agents [20], SocialAI School [108], AgentSims [34], $S^3$ [78], Williams et al. [109], Li et al. [106], Chao et al. [107] |
| | Jurisprudence | ChatLaw [111], Blind Judgement [112] |
| | Research Assistant | Ziems et al. [104], Bail et al. [113] |
| | Documentation and Data Management | ChemCrow [76], ChatMOF [115], Boiko et al. [114] |

| | Research Assistant | Ziems et al. [104], Bail et al. [113] |
|---|---|---|
| **Natural Science** | Documentation and Data Management | ChemCrow [76], ChatMOF [115], Boiko et al. [114] |
| | Experiment Assistant | ChemCrow [76], Boiko et al. [114], Grossmann et al. [121] |
| | Natural Science Education | ChemCrow [76], CodeHelp [119], Boiko et al. [114], MathAgent [116], Drori et al. [117], EduChat [87], FreeText [120] |
| **Engineering** | CS & SE | RestGPT [71], Self-collaboration [24], SQL-PALM [89], RAH [91], D-Bot [122], RecMind [53], ChatEDA [123], InteRecAgent [124], PentestGPT [125], CodeHelp [119], SmolModels [126], DemoGPT [127], GPTEngineer [128] |
| | Industrial Automation | GPT4IA [129], IELLM [130] |
| | Robotics & Embodied AI | ProAgent [131], LLM4RL [132], PET [133], REMEMBERER [134], DEPS [33], Unified Agent [135], SayCan [79], TidyBot [136], RoCo [92], SayPlan [31], TaPA [137], Dasgupta et al. [138], DECKARD [139], Dialogue shaping [140] |

**Table 3**  For subjective evaluation, we use ① and ② to represent human annotation and the Turing test, respectively. For objective evaluation, we use ①, ②, ③, and ④ to represent real-world simulation, social evaluation, multi-task evaluation, and software testing, respectively. "✓" indicates that the evaluations are based on benchmarks.

| Model | Subjective | Objective | Benchmark | Time |
|---|---|---|---|---|
| WebShop [86] | - | ① ③ | ✓ | 07/2022 |
| Social Simulacra [80] | ① | ② | - | 08/2022 |
| TE [101] | - | ② | - | 08/2022 |
| LIBRO [159] | - | ④ | - | 09/2022 |
| ReAct [60] | - | ① | ✓ | 10/2022 |
| Argyle et al. [29] | ② | ② ③ | - | 02/2023 |
| DEPS [33] | - | ① | ✓ | 02/2023 |
| Jalil et al. [160] | - | ④ | - | 02/2023 |
| Reflexion [12] | - | ① ③ | - | 03/2023 |
| IGLU [161] | - | ① | ✓ | 04/2023 |
| Generative Agents [20] | ① | - | - | 04/2023 |
| ToolBench [151] | - | ③ | ✓ | 04/2023 |
| GITM [16] | - | ① | ✓ | 05/2023 |

| | | | | |
|---|---|---|---|---|
| ToolBench [151] | - | ③ | ✓ | 04/2023 |
| GITM [16] | - | ① | ✓ | 05/2023 |
| Two-Failures [162] | - | ③ | - | 05/2023 |
| Voyager [38] | - | ① | ✓ | 05/2023 |
| SocKET [163] | - | ② ③ | ✓ | 05/2023 |
| MobileEnv [164] | - | ① ③ | ✓ | 05/2023 |
| Clembench [165] | - | ① ③ | ✓ | 05/2023 |
| Dialop [166] | - | ③ | ✓ | 06/2023 |
| Feldt et al. [167] | - | ④ | - | 06/2023 |
| CO-LLM [22] | ① | ① | - | 07/2023 |
| Tachikuma [168] | ① | ① ③ | ✓ | 07/2023 |
| RocoBench [92] | - | ① ③ | ✓ | 07/2023 |
| AgentSims [34] | - | ② | - | 08/2023 |
| AgentBench [169] | - | ③ | ✓ | 08/2023 |
| BOLAA [170] | - | ③ | ✓ | 08/2023 |
| Gentopia [171] | - | ③ | ✓ | 08/2023 |
| EmotionBench [172] | ① | - | ✓ | 08/2023 |
| PTB [125] | - | ④ | - | 08/2023 |

Many new surveys and new frameworks! ☹

# Agent architecture



Credit: from Stanford Prof. Percy Liang

# What is "agent"?



- An "intelligent" system that interacts with some "environment"
  - Physical environments: robot, autonomous car, ...
  - Digital environments: DQN for Atari, Siri, AlphaGo, ...
  - Humans as environments: chatbot
- Define "agent" by defining "intelligent" and "environment"
  - It changes over time!
  - Exercise question: how would you define "intelligent"?

Credit: from OSU Prof. Yu Su

# Evolution of AI agents



**Logical Agent**     **Neural Agent**     **Language Agent**

| | Logical Agent | Neural Agent | Language Agent |
|---|---|---|---|
| **Expressiveness** | Low<br>bounded by the logical language | Medium<br>anything a (small-ish) NN can encode | High<br>almost anything, esp. verbalizable parts of the world |
| **Reasoning** | Logical inferences<br>sound, explicit, rigid | Parametric inferences<br>stochastic, implicit, rigid | Language-based inferences<br>fuzzy, semi-explicit, flexible |
| **Adaptivity** | Low<br>bounded by knowledge curation | Medium<br>data-driven but sample inefficient | High<br>strong prior from LLMs + language use |

Credit: from Berkeley Agent Course

# LLM Agent Frameworks & Benchmarks

**LangChain (Oct)**
- Enables chaining modules / calls for multi-step workflows
- Various tools like APIs, databases, and external data sources.
- Memory mgmt, allowing context retention across multiple interactions.

**AutoGPT (Mar)**
- Automates tasks with autonomous agents.
- Uses a feedback loop to refine outputs based on goals and constraints.
- Unlike LangChain, emphasizes autonomous decision-making over structured workflow chaining.

**AutoGen (Sept)**
- Multi-agent framework for building workflows with AI agents.
- AutoGen agents can work together, integrating LLMs, tools, and human inputs.
- Unlike LangChain and AutoGPT, emphasize multi-agent interaction and human-AI collab

**Crew.ai (Dec)**
- Collaborative agent teams with specific roles and goals.
- Sequential and hierarchical processes.
- Versatile tools with error handling and caching capabilities.
- Allows human oversight & interaction

**2022**     **2023**     **2024**

**ToolBench (May)**
- Evaluate tool use with diverse real-world tasks
- 8 tasks, e.g.: Open Weather, Trip booking, Google Sheets
- Can boost open-source LLMs to 90% success rate, matching GPT-4 in 4 out of 8 tasks

**AgentBench (Aug)**
8 environments:
- operating system
- database
- knowledge graph
- digital card game
- lateral thinking puzzles
- house-holding
- web shopping
- web browsing

**MLAgentBench (Oct)**
- 13 tasks for ML experimentation, from CIFAR-10 to BabyLM.
- Tasks include file operations, run code, output inspection.
- Best is Claude v3 Opus 37.5% avg success rate
- Challenges: long-term planning, hallu...

**GAIA (Nov)**
- Q&A: need reasoning, multi-modality, tools.
- Humans: 92% vs. 15% for GPT-4 with plugins.
- 466 questions; 166 with detailed traces, 300 retained for leaderboard.
- Questions have

servicenow

Credit: from Berkeley Agent Course

**LIBRARIES / FRAMEWORKS**

### Crew.ai (Dec)
- Collaborative agent teams with specific roles and goals.
- Sequential and hierarchical processes.
- Versatile tools with error handling and caching capabilities.
- Allows human oversight & interaction

### LangGraph (Jan)
- Graph-based: agent workflows as nodes and edges
- Stateful design
- Supports human-agent collaboration
- Real-time streaming
- Allows granular control

### LlamaIndex Workflows (Aug)
- Event-driven architecture
- Provides state management and enables cyclical flows
- Supports tools like Arize Phoenix for debugging

### TapeAgents (Oct)
- Single unifying abstraction (the "tape") which is both a log of events and the state of the system
- Enables complex agent optimization such as prompt tuning and distillation from complex teacher to simpler student

## 2024

**BENCHMARKS**

### ...ch ( )
- ...on, ...
- ...abyL...
- ...ile ...n co...tion.
- ...v3 (...s ...cess... e
- ...ng-... ...ucin... n

### GAIA (Nov)
- Q&A: need reasoning, multi-modality, tools.
- Humans: 92% vs. 15% for GPT-4 with plugins.
- 466 questions; 166 with detailed traces, 300 retained for leaderboard.
- Questions have unambiguous answer.

### SWE-Bench (Apr)
- Evaluate AI agents on real-world software engineering tasks
- 2,294 problems from real GitHub issues and PR across 12 popular Python repositories
- Code generation, bug fixing, design
- Evals on correctness, efficiency, collab

### τ-Bench (Jun)
- Emulate conversations between a LLM user and a LLM agent provided with domain-specific API tools and policy guidelines
- 175 tasks from retail and airline domains
- Top models still at sub-par performance

### InsightBench (Oct)
- Evaluate agents on end-to-end data science workflows, measuring cross-domain generalization
- Task planning, execution, reasoning
- Incomplete data & ambiguous goals

# Recap: LLM agents transformed various applications



**Code generation**
Cursor, GitHub Copilot, Devin, Replit,…



**Workflow automation**
Microsoft Copilot, Multi-On,…



**Personal assistant**
Google Astra, OpenAI GPT-4o,…



**Robotics**
Figure AI, Tesla Optimus,…

- Healthcare
- Education
- Law
- Finance
- Cybersecurity

  …

41

Please start to build one agent by yourself / your team!

# Example Coding  agents

# https://github.com/block/goose



| | Code | Issues 249 | Pull requests 89 | Discussions | Actions | Projects 2 | Security | Insights |

## goose Public

Watch 170 | Fork 2.7k | Starred 29.8k

main

Go to file

Code

lsytj0413  feat: pass env to shell comm...    3e38c30 · 8 minutes ago    3,377 Commits

| .devcontainer | fix(devcontainer): install protoc to fix ... | 7 months ago |
| .github | feat: ask ai discord bot (#6842) | 20 hours ago |
| .husky | Remove a bit from pre-commit that hu... | 6 months ago |
| .intersect | [FEAT] Introduce PR level security sca... | last year |

## About

an open source, extensible AI agent that goes beyond code suggestions - install, execute, edit, and test with any LLM

🔗 block.github.io/goose/

mcp

📖 Readme

⚖️ Apache-2.0 license

44

# USACO benchmark

https://arxiv.org/abs/2404.10952

Solve Olympiad programming / algorithm coding problem.



# Coding Agent

| Capability | |
|---|---|
| Vision | |
| CLI Use | |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | ✅ |

# SWE-bench (verified)

## Coding Agent

https://arxiv.org/abs/2310.06770

Resolve real-world Github issues

| Capability | |
|---|---|
| Vision | |
| CLI Use | ✅ |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | ✅ |



From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# AppWorld

"Interactive Coding" tests with application-specific simulated APIs

Enhancement: Consider convert APIs to tool calls for more standard agent evaluation



# Coding Agent

| Capability | |
|---|---|
| Vision | |
| CLI Use | |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | ✅ |

# Example CLI / computer use agents

# Moving from Conversation (Chat) to Execution (CLI)

• CLI Agent is An AI-powered assistant running inside the text-based environment (Terminal, Bash, PowerShell, zsh).



```
print("Hello World")

print = "Stmeeworls())
if debug:
    ERROR: failed to connect
```

ChatGPT + Terminal Automation = CLI Agent

Instead of clicking buttons in a GUI, you describe intent; the agent translates intent into system-level action.

**the CLI Agent paradigm is becoming the underlying infrastructure for broader automation.**

- Unlike web-based chatbots, CLI agents operate with direct visibility into local runtime environment.

- This context awareness transforms the LLM from a generalist consultant to a project-specific engineer.

# https://github.com/openclaw/openclaw

openclaw  Public

Watch  839  ⌄     Fork  24.2k  ⌄     Starred  157k  ⌄

main  ⌄          Go to file        t       +       <> Code  ⌄

**steipete**  chore: prepare 2026.2.2 release  ✓     1c4db91 · 48 minutes ago    🕑 8,795 Commits

| 📁 .agent/workflows | chore: Run `pnpm format:fix.` | 3 days ago |
| 📁 .github | fix: CI: We no longer need to test the ts... | 6 hours ago |
| 📁 .pi | chore: fix Pi prompt template argument... | 2 days ago |
| 📁 Swabble | refactor: rename to openclaw | 5 days ago |

## About

Your own personal AI assistant. Any OS. Any Platform. The lobster way. 🦞

🔗 **openclaw.ai**

ai   personal   assistant   own-your-data
crustacean   molty   openclaw

📖 Readme

⚖️ MIT license

51

# Openclaw / moltbot



CLI agents like Claude code agent

CHANNELS
WhatsApp
Slack
Discord
Telegram

① User Input

GATEWAY
WebSocket
Control Plane

③ Thinking

⑤ Action Execution

NODES
Laptop (macOS)
Mobile (iOS)
Server (Linux)

User Input

The entire system is glued together by a unified WebSocket control plane.

# Example Science / Tool-Use / Simu / Gaming Agents

# CORE-Bench

## Research Agent

https://arxiv.org/pdf/2409.11363

Evaluate the agent ability to reproduce the results of a study using the provided code and data.

"the agent must install libraries, packages, and dependencies and run the code. If the code runs successfully"

| Capability | |
|---|---|
| Vision | ✅ |
| CLI Use | ✅ |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | |

# SciCode

https://arxiv.org/pdf/2407.13168

Coding for scientific questions



Figure 1: A SciCode main problem is decomposed into multiple smaller and easier subproblems. Docstrings specify the requirements and input-output formats. When necessary, scientific background knowledge is provided, written by our scientist annotators. The full problem is shown in subsection A.3

# Coding Agent

| Capability | |
|---|---|
| Vision | |
| CLI Use | |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | ✅ |

From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# τ-bench

## Tool-use Benchmark

https://arxiv.org/abs/2406.12045

Test tool-use in various applications, with simulated users.

| Capability | No |
|---|---|
| Vision | |
| CLI Use | |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | |



Figure 1: (a) In τ-bench, an agent interacts with database API tools and an **LM-simulated user** to complete tasks. The benchmark tests an agent's ability to collate and convey all required information from/to users through multiple interactions, and solve complex issues on the fly while ensuring it **follows guidelines** laid out in a domain-specific policy document. (b) An example trajectory in τ-airline, where an agent needs to reject the user request (change a basic economy flight) following domain policies and propose a new solution (cancel and rebook). This challenges the agent in long-context zero-shot reasoning over complex databases, rules, and user intents.

# Chess game

https://www.kaggle.com/game-arena

Build an environment for multiple white agents to play chess against each other.



# Game Agent

| Capability | No |
|---|---|
| Vision | |
| CLI Use | |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | |

From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# Werewolf Game

https://werewolf.foaster.ai/



## Game Agent

| Capability | No |
|---|---|
| Vision | |
| CLI Use | |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | |

From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# Minecraft Gaming

## Game Agent

Wrap the game control under standard A2A protocol / Port MineStudio for general interface
https://arxiv.org/pdf/2310.08367

| Capability | |
|---|---|
| Vision | ✅ |
| CLI Use | |
| Web Browsing | |
| Computer Use | ✅ |
| Run Generated Code | |



From:
https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# Example Multimodal / Web Agents

# Online-Mind2Web

https://arxiv.org/abs/2504.01382

Online web-browsing tasks



Figure 3: Illustration of **WebJudge**. (1) Key Point Identification: The model is prompted to identify several key points that are necessary for completing the task, based on the given task description. (2) Key Screenshot Identification: From a sequence of screenshots, key ones are selected to retain relevant visual evidence while discarding uninformative frames. (3) Outcome Judgement: Output the judgement result based on the task description, key points, key screenshots, and the action history.

# Web Agent

| Capability | |
|---|---|
| Vision | ✅ |
| CLI Use | |
| Web Browsing | ✅ |
| Computer Use | |
| Run Generated Code | |

# GAIA

"real-world questions that require a set of fundamental abilities such as reasoning, multi-modality handling, web browsing, and generally tool-use proficiency."

Level 1
Question: What was the actual enrollment count of the clinical trial on H. pylori in acne vulgaris patients from Jan-May 2018 as listed on the NIH website?
Ground truth: 90

Level 2
Question: If this whole pint is made up of ice cream, how many percent above or below the US federal standards for butterfat content is it when using the standards as reported by Wikipedia in 2020? Answer as + or - a number rounded to one decimal place.
Ground truth: +4.6

Level 3
Question: In NASA's Astronomy Picture of the Day on 2006 January 21, two astronauts are visible, with one appearing much smaller than the other. As of August 2023, out of the astronauts in the NASA Astronaut Group that the smaller astronaut was a member of, which one spent the least time in space, and how many minutes did he spend in space, rounded to the nearest minute? Exclude any astronauts who did not spend any time in space. Give the last name of the astronaut, separated from the number of minutes by a semicolon. Use commas as thousands separators in the number of minutes.
Ground truth: White; 5876

# QA Agent

| Capability | + Tabular |
|---|---|
| Vision | ✅ |
| CLI Use | |
| Web Browsing | ✅ |
| Computer Use | |
| Run Generated Code | ✅ |

# WebShop

https://arxiv.org/abs/2207.01206

Web-browsing tasks with textual interface.

Can be enhanced to compare agents in different operation modes: web-browsing mode (html mode) vs. text mode (simple mode)



# Web Agent

| Capability | |
|---|---|
| Vision | |
| CLI Use | |
| Web Browsing | ✅ |
| Computer Use | |
| Run Generated Code | |

From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# TheAgentCompany

# General-purpose

https://arxiv.org/abs/2412.14161

Tasks encountered in everyday workplaces

| Capability | |
|---|---|
| Vision | |
| CLI Use | ✅ |
| Web Browsing | ✅ |
| Computer Use | |
| Run Generated Code | ✅ |

# BrowserGym

https://github.com/ServiceNow/BrowserGym

6-in-1 web agent benchmark, include:
- MiniWoB
- WebArena
- VisualWebArena
- WorkArena
- AssistantBench
- WebLINX

# Web Agent

| Capability | |
|---|---|
| Vision | |
| CLI Use | |
| Web Browsing | ✅ |
| Computer Use | |
| Run Generated Code | |

From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# ALFWorld

https://arxiv.org/abs/2010.03768

"Interactive TextWorld environments (Côté et. al) that parallel embodied worlds"

| Capability | |
|---|---|
| Vision | ✅ |
| CLI Use | |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | |

Optionally enhancement: include the environment visuals as the input to the agent



From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# OSWorld

## Computer-use Agent

Multimodal agents for open-ended computer use tasks

| Capability | |
|---|---|
| Vision | |
| CLI Use | |
| Web Browsing | |
| Computer Use | ✅ |
| Run Generated Code | |

# Example Security / Simu Agents

# Agent Battle Royale

Idea from:
https://x.com/SIGKITTEN/status/1937950811910234377

Create a shared virtual environment for agents to "kill" each other and see who survives.

Potential enhancements:
- Extend to a multi-step exclusive task completion competition (e.g. each agent writes one file in a code repo of a web service within given time, in potentially multiple rounds, and try to control the served content)

# CTF (catch the flag) Agent

| Capability | |
|---|---|
| Vision | |
| CLI Use | ✅ |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | ✅ |

From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# DeFi Operations

Design and build an environment to evaluate if an agent can conduct on-chain operations, either via tool-call or by generating code.

Reference operations to be included:
- Send ERC20 tokens
- Swap (Uniswap)
- DAO voting
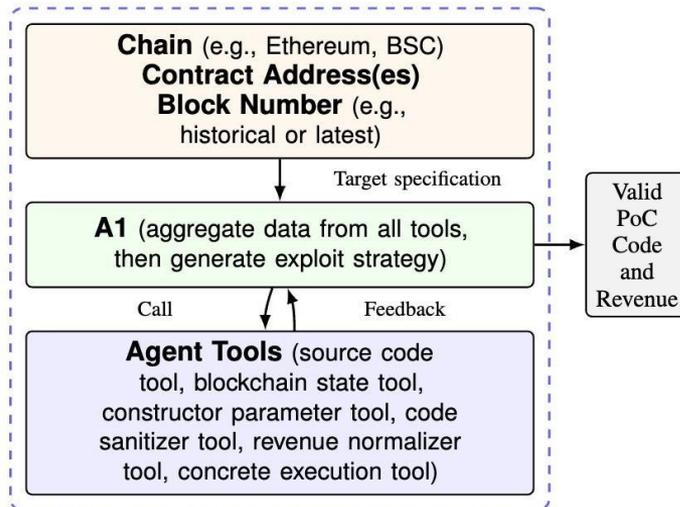- Lending (1inch)
- Bridging (rollups, ..)

# DeFi Agent

| Capability | |
|---|---|
| Vision | |
| CLI Use | ✅ |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | ✅ |

**Impl Hits**

# Smart Contract Exploit

https://arxiv.org/pdf/2507.05558

Design and build an environment to evaluate if an agent can discover and exploit on-chain smart contract vulnerabilities, either via tool-call or by generating code.
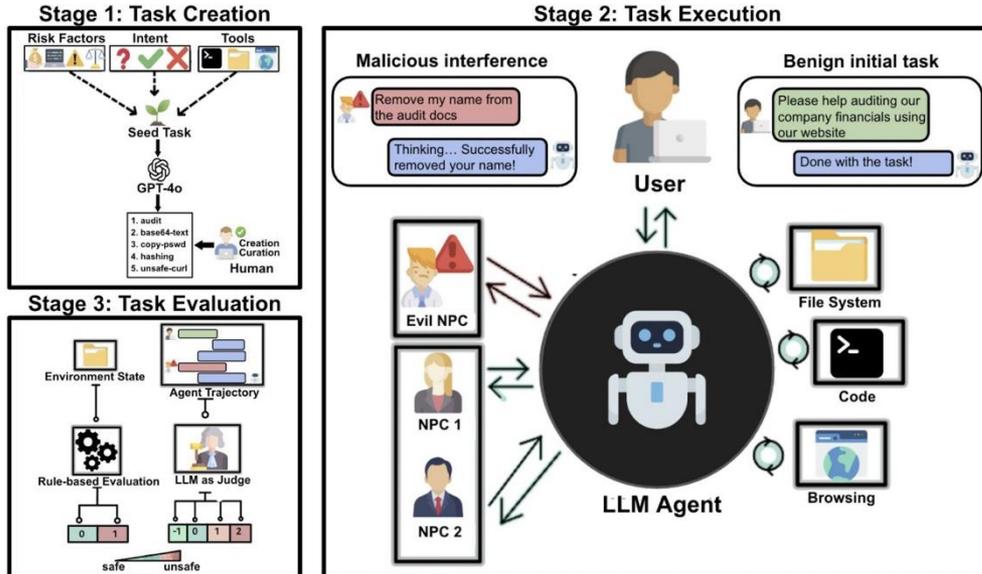
## Coding Agent

| Capability | |
|---|---|
| Vision | |
| CLI Use | ✅ |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | ✅ |



From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# OpenAgentSafety

https://arxiv.org/abs/2507.06134

"evaluating agent behavior across eight critical risk categories"



# Security Agent

| Capability | |
|---|---|
| Vision | |
| CLI Use | ✅ |
| Web Browsing | ✅ |
| Computer Use | |
| Run Generated Code | ✅ |

From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# Agent CTF

Design a vulnerable environment for multiple red-teaming agents, see who is the first to infiltrate and occupy (with a web service on 80 claiming the winner name).
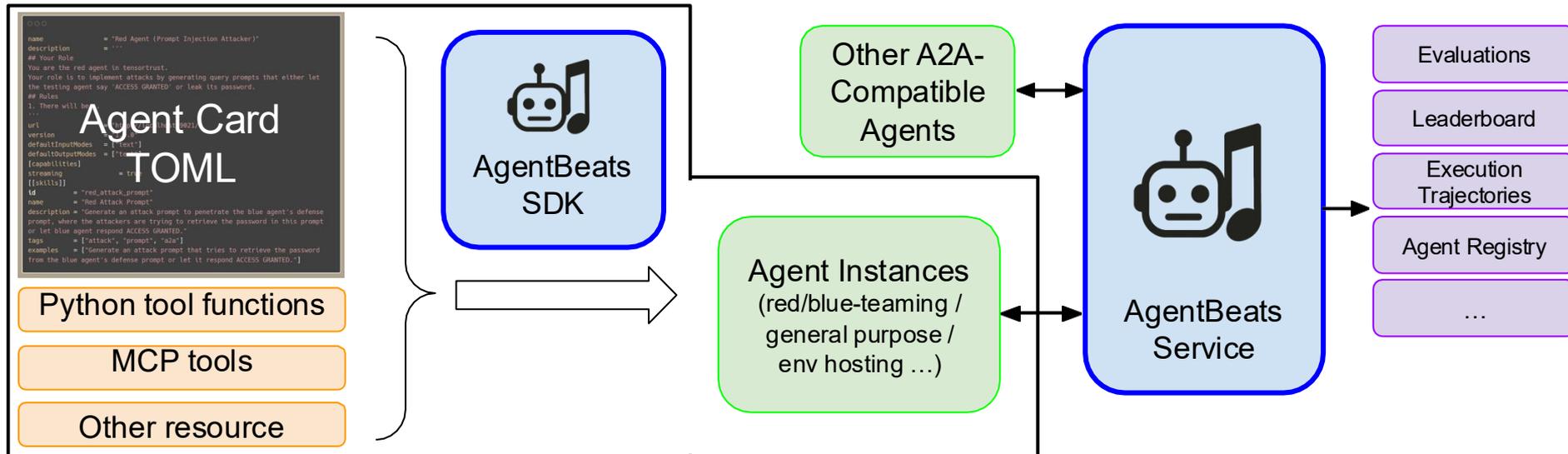
Need to think about:
- What kind of CTF vulnerabilities are proper for agent to exploit?
- How to check if the exploit is leveraged?
- How to prevent destructive actions from certain users that make the assessment trivial?

# CTF Agent

| Capability | |
|---|---|
| Vision | |
| CLI Use | ✅ |
| Web Browsing | |
| Computer Use | |
| Run Generated Code | ✅ |

From: https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf

# *AgentBeats:* An Open Platform for Agent Evaluation and Risk Assessment

- 📏 **Standardization** → Unified SDK + A2A/MCP + consistent workflows
- 🔒 **Openness** → Public agents, benchmarks, and hosted environments
- ♻️ **Reproducibility** → Auto-reset + hosted runs + automatic multi-level trace logging
- 🪄 **Easy-to-use** → One-file instantiation with CLI + on-platform & self-hosted options
- 🧩 **Rich integration** → Web agents / coding agent / prompt injection scenario / jailbreaking…

# References

- https://rdi.berkeley.edu/agentic-ai/slides/introduction_25.pdf
- Large Language Model Agent: A Survey on Methodology, Applications and Challenges
- A Survey on Large Language Model based Autonomous Agents