



# *AutoAgents: A Framework for Automatic Agent Generation*

Künstliche Intelligenz Seminar 2025

Prof. Dr. Schmidt-Schuß

Mahmoud Alhag Ali





## **Autoren:**

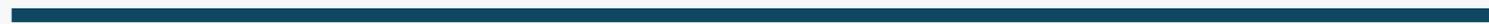
Guangyao Chen, Siwei Dong,  
Yu Shu, Ge Zhang,  
Jaward Sesay, Börje Karlsson,  
Jie Fu and Yemin Shi

## **Instituten**

Peking University,  
Hong Kong University of Science and Technology  
Beijing Academy of Artificial Intelligence  
University of Waterloo  
{gy.chen, shiyemin}@pku.edu.cn, jiefu@ust.hk

## **Conferenc**

Proceedings of the Thirty-Third International Joint Conference on  
Artificial Intelligence (IJCAI-24), site (22 - 30)



# Agenda

1. Definitionen
2. Vorarbeiten
3. Entwurfsphase
4. Agent-Architektur und Komponentendesign
5. Plan Erstellung
6. Ausführungsphase
7. Workflow
8. Experimentelle Validierung und Leistungsanalyse
9. Praktische Anwendung und Fallstudie

# Einleitung



---

Die meisten bestehenden LLM-basierten Multi-Agent-Ansätze basieren auf vordefinierten Agenten zur Bewältigung einfacher Aufgaben, was die Anpassungsfähigkeit der Multi-Agent-Kollaboration an verschiedene Szenarien stark einschränkt. AutoAgents löst dieses fundamentale Problem durch die adaptive Generierung und Koordination mehrerer spezialisierter Agenten zur Bildung eines KI-Teams entsprechend unterschiedlicher Aufgaben

---





# 1. Definitionen

## **AI Agent:**

Ein AI Agent (deutsch: KI-Agent) ist ein autonomes Softwaresystem, das mithilfe künstlicher Intelligenz in der Lage ist, eigenständig Aufgaben in einer definierten Umgebung auszuführen, um bestimmte Ziele zu erreichen.

## **LLM Agent**

LLM-basierte autonome Agenten: LLMs werden häufig als Kernsteuerungen für autonome Agenten verwendet, die bestimmte Ziele erreichen können.

## **Multiagentensystem:**

Ein Multiagentensystem (kurz: MAS) bezeichnet ein System, das aus mehreren autonomen Agenten besteht, die miteinander und mit ihrer Umgebung interagieren, um individuelle oder gemeinsame Ziele zu erreichen.



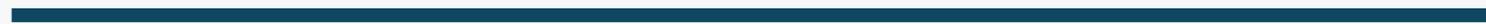


## 2. Vorarbeiten

**Auto-GPT [Gravitas, 2023]:** ist eine frühe Arbeit, die LLM als KI-Agent einsetzt, jedoch unterstützt es keine Multi-Agenten-Kolaboration und kann nur isoliert arbeiten.

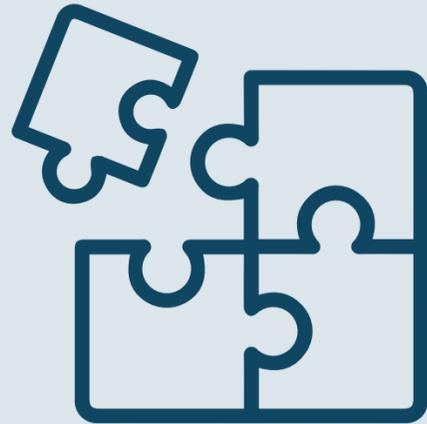
**SSP [Wang et al., 2023c]:** hat die Aufgabenlösungsfähigkeiten von LLMs durch Integration von Multi-Agenten-Diskussionssystem. Es erlaubt LLMs für oder gegen eine Aufgabe zu argumentieren.

**AgentVerse [Chen et al., 2023a]:** Framework für die automatische Erzeugung unbegrenzter Agenten, es generiert den Ausführungsplan durch die Diskussionen der generierten Agenten und fügt Bewertungsstrategien für die zyklische Ausführung hinzu



## 3. Zweistufige Architektur

---



### Entwurfsphase

- Synthetisiert ein angepasstes Agententeam und Ausführungsplan
- Planner, Agent Observer, Plan Observer



### Ausführungsphase

- Verfeinerung den Plan und produziert das endgültige Ergebnis
- Action Observer

# 3.1 Entwurfsphase



synthetisiert ein angepasstes Agententeam und einen Ausführungsplan durch kollaborative Diskussion zwischen drei vordefinierten Agenten

## Planner (P)

Generiert und verfeinert ein Agententeam sowie einen Ausführungsplan basierend auf dem Aufgabeninhalt

---

## Agent Observer (O\_agent)

Bietet Vorschläge zur Rationalität der Agententeammitglieder und deren Übereinstimmungsgrad mit der Aufgabe

---



## Plan Observer (O\_plan)

Liefert Vorschläge zur Rationalität des Ausführungsplans und dessen Übereinstimmungsgrad mit der Aufgabe und dem Agententeam

---



## 3.2 Ausführungsphase



### Ausführungsphase

Die Ausführungsphase verfeinert den Plan durch Inter-Agent-Kollaboration und Feedback und produziert das endgültige Ergebnis

---



### Action Observer (O\_action)

fungiert als Aufgabenmanager für die verschiedenen Agenten, weist ihnen unterschiedliche Aufgaben zu, überprüft die Ausführungsergebnisse jedes Agenten und passt den Ausführungsplan dynamisch basierend auf dem Ausführungsstatus an

---



# 4. Agentenstruktur und Komponentendesign

## Prompt (P)

- Bietet eine detaillierte und angepasste Darstellung der Expertenidentität für jeden spezifischen Agenten, bestehend aus Profil, Ziel und Beschränkungen

## Description (D)

- Verleiht zusätzliche konkrete Identität zur Etablierung einer umfassenderen Rolle

## Toolset (T)

- Stattet den Agenten mit ausgewählten Tools aus einem vordefinierten Set aus, um Entscheidungsverwirrung durch übermäßige Tools zu verhindern

## Suggestions (S)

- Liefert Vorschläge für jeden Agenten zur Ausführung der aktuellen Aufgabe

# 4. Agent Observer

## Planner (P)

- Auf der Grundlage der vom Planner wird Agentenliste  $\{A_1, A_2, \dots, A_n\}$  erstellt

## Agent Observer (O\_agent)

1. Prüft, ob jeder Agent den genannten Spezifikationen entspricht und identifiziert alle fehlenden Elemente  $\{P, D, T\}$
2. bewertet die Kompatibilität jedes Agenten mit der Aufgabe anhand seiner Beschreibungsinformationen und des Aufgabeninhalts.
3. Prüft die Agentenliste auf überflüssige oder fehlende Rollen und eliminiert oder ergänzt sie entsprechend.

# 5. Plan Erstellung

## Zeitgleiche Ausführungsplanung:

- Während Agenten für eine Aufgabe erstellt werden, entwickelt der Planer parallel einen Ausführungsplan zur Bewältigung der Aufgabe.

## Planerstellung:

- Der Planer legt für die Aufgabe konkrete Einzelschritte  $\{S_1, S_2, \dots, S_n\}$  fest.
- Jeder Schritt enthält:
  - Den zuständigen Agenten ( $A(j)$ ),
  - die notwendigen Eingaben und
  - die erwarteten Ergebnisse.

## Planprüfung durch den Plan Observer:

- Der Plan Observer kontrolliert, ob jedem Schritt ein klarer Agent zugeordnet ist und ob die Inhalte der Einzelschritte verständlich und konsistent sind.
- Außerdem prüft er, ob die Schrittfolge vollständig und geeignet ist, um die Aufgabe vollständig zu lösen.

# 5. Plan Erstellung

## **Iterative Verbesserung:**

- Durch wiederholten Austausch zwischen Planer und Plan\_Observer wird der Plan schrittweise optimiert.

## **Abschluss:**

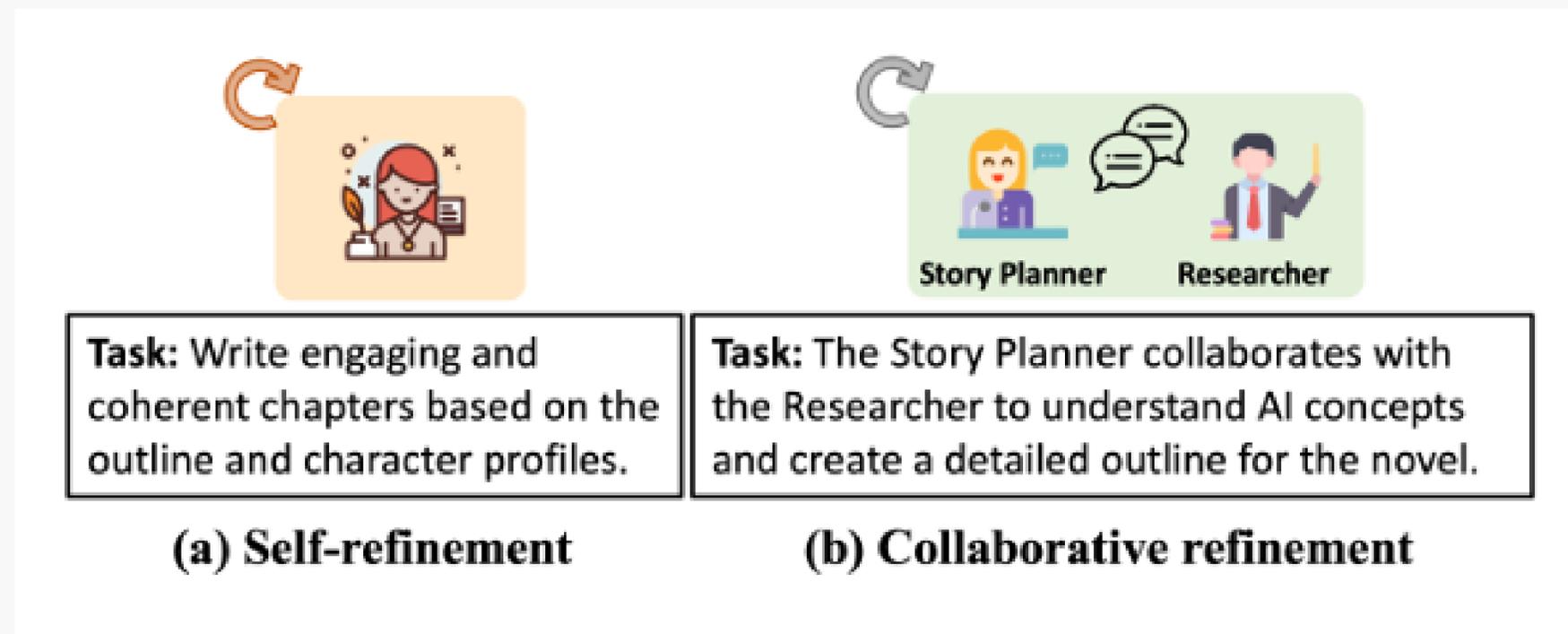
- Nach mehreren Überprüfungsschleifen entsteht ein konsistenter, vollständiger Ausführungsplan zur Erreichung der vorgegebenen Aufgabe.

# 5.1 Aktionen zur Aufgabenausführung



Der Ausführungsplan umfasst zwei Aktionen der Aufgabenausführung:

- Selbstverfeinerung durch einen einzelnen Agenten
- Kollaborative Verfeinerung durch mehrere Agenten



## 6. Ausführungsphase



**Action Observer O\_action:** fungiert als Aufgabenmanager für die verschiedenen Agenten, indem er ihnen verschiedene Aufgaben zuweist, die Ausführungsergebnisse der einzelnen Agenten überprüft und den Ausführungsplan auf der Grundlage des Ausführungsstatus dynamisch anpasst.

Dieser Mechanismus der Verfeinerung und Kommunikation wiederholt sich so lange, bis der Aktion\_Observer eine einstimmige Einigung über die Ausführungsantworten erzielt oder der Prozess seine maximale Iteration erreicht.





## 6.1 Mechanismus zum Wissensaustausch.

Wenn jedoch die Anzahl der Agenten groß ist und ein einzelner Agent mehr Selbstiterationen durchführt, generiert er mehr historische Informationen. Aufgrund der Token-Beschränkung von LLM-Modellen können diese oft nicht alle Informationen umfassen. In diesem Rahmen gibt es also

- ein Kurzzeitgedächtnis,
- ein Langzeitgedächtnis,
- ein dynamisches Gedächtnis.



# 6.1 Innovatives Speichersystem

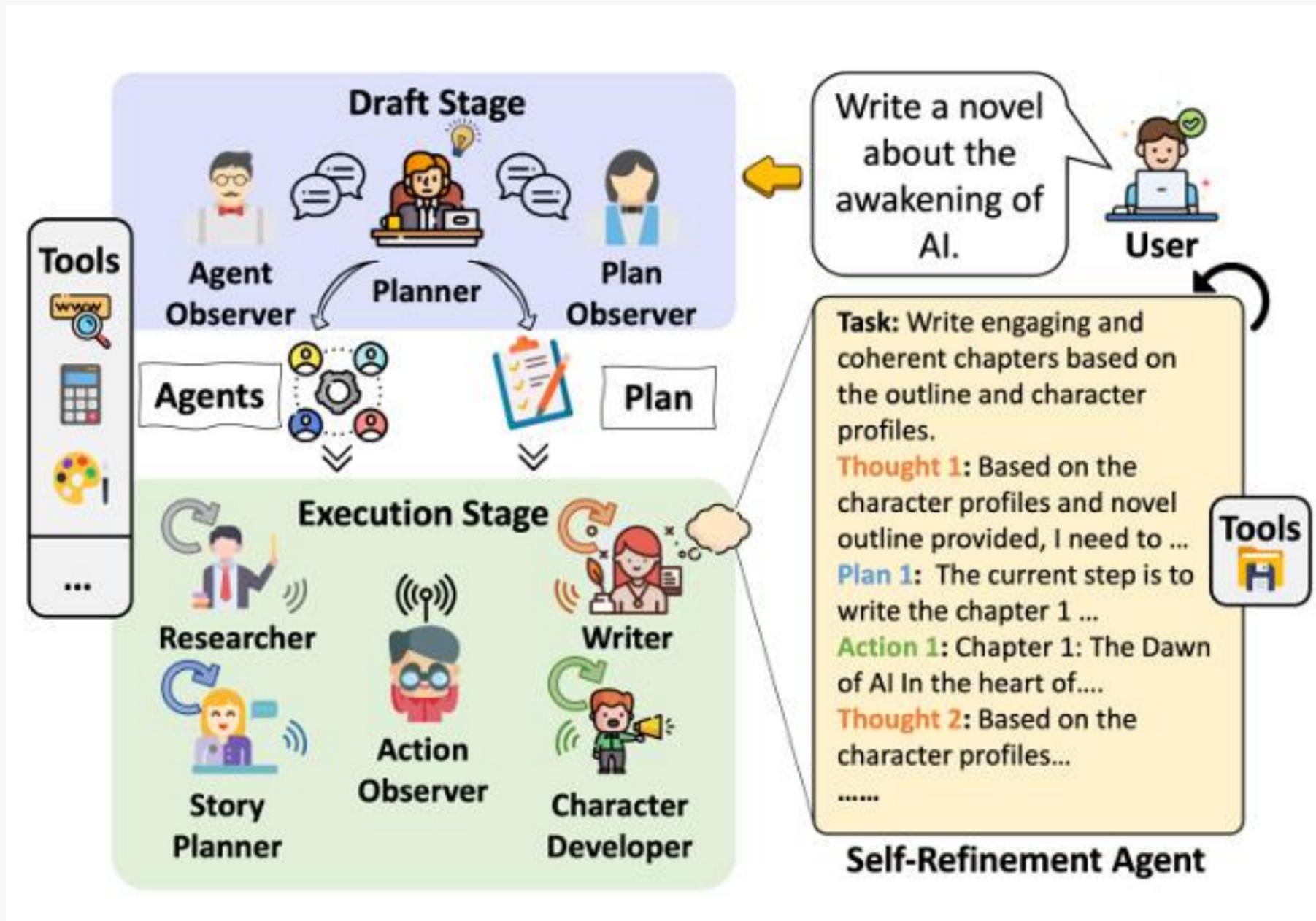


AutoAgents implementiert ein dreischichtiges Speichersystem zur Bewältigung komplexer Kontextdaten

Speichertyp	Zweck	Inhalt
Kurzzeitgedächtnis	Einzelne Aktion	Zwischengedanken, Strategien und Ergebnisse während Self-/Collaborative Refinement
Langzeitgedächtnis	Historische Trajektorie	Ausgeführte Ergebnisse jeder Aufgabe und Synthese wichtiger Feedback-Informationen
Dynamisches Gedächtnis	Spezialisierte Aufmerksamkeit	Vom Action Observer extrahierte Hilfsinformationen für spezifische Aktionsanforderungen



# 7. Workflow



## Algorithm 1 AutoAgents Execution Process.

**Input:** User task/Question

**Output:** Task solution/Answer

- 1: **Drafting Stage**
- 2: Initialize Planner  $\mathcal{P}$ , Agent Observer  $\mathcal{O}_{\text{agent}}$ , and Plan Observer  $\mathcal{O}_{\text{plan}}$ .
- 3:  $\mathcal{P}$  generates initial agent team  $\{\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n\}$  and execution plan  $P = \{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_n\}$ .
- 4: **repeat**
- 5:  $\mathcal{O}_{\text{agent}}$  provides feedback on agent team.
- 6:  $\mathcal{P}$  refines agent team based on feedback.
- 7:  $\mathcal{O}_{\text{plan}}$  provides feedback on execution plan.
- 8:  $P$  refines execution plan based on feedback.
- 9: **until** No feedback or reached the maximum iteration limit.
- 10: **Execution Stage:**
- 11: Initialize Action Observer  $\mathcal{O}_{\text{action}}$  and long-term memory  $M_L$ .
- 12: **for**  $\{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_n\}$  **do**
- 13:  $\mathcal{O}_{\text{action}}$  generate dynamic memory  $M_D$ .
- 14:  $\mathcal{O}_{\text{action}}$  assign task  $\mathcal{S}_k$  and  $M_D$  to corresponding agents  $\{\mathcal{A}_i, \dots, \mathcal{A}_j\}$ .
- 15: Initialize short-term memory  $M_S$ .
- 16: **repeat**
- 17: **for**  $\{\mathcal{A}_i, \dots, \mathcal{A}_j\}$  **do**
- 18: Agent  $\mathcal{A}_m$  analyzes  $\mathcal{S}_k$ ,  $M_S$  and  $M_D$ .
- 19: Agent  $\mathcal{A}_m$  plans the current step and executes this step.
- 20: The execution result is stored in  $M_S$ .
- 21: **end for**
- 22: **until** No step or reached the maximum iteration limit.
- 23: The execution results of task  $\mathcal{S}_k$  are stored in  $M_L$ .
- 24:  $\mathcal{O}_{\text{action}}$  coordinates  $\{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_n\}$  and monitors execution.
- 25: **end for**
- 26: **return** Execution results of final step.

# 8. Experimentelle Validierung und Leistungsanalyse



Open-ended Question Answering

AutoAgents wurde auf dem MT-bench-Datensatz mit 80 hochwertigen offenen Fragen aus verschiedenen Kategorien wie gesunder Menschenverstand, kontrafaktische Szenarien und Codierung getestet

<b>Evaluator</b>	<b>v.s. GPT-4</b>	<b>v.s. AgentVerse</b>
FairEval [Wang <i>et al.</i> , 2023b]	76.3%	91.3%
HumanEval	62.5%	77.5%



# 8. Experimentelle Validierung und Leistungsanalyse



## Trivia Creative Writing

Diese Aufgabe fordert die Fähigkeiten großer Sprachmodelle heraus, vielfältige Informationen aus ihrem internen selbst komprimierten Wissen abzurufen und zu integrieren. Das Modell muss eine kohärente Geschichte um ein gegebenes Thema verfassen und dabei Antworten auf N Trivia-Fragen einbeziehen.

Methods	N (# trivia questions) = 5		N (# trivia questions ) = 10	
	Score (%)	$\Delta$ (v.s Standard %)	Score (%)	$\Delta$ (v.s Standard %)
Standard	74.6	0.0%	77.0	0.0%
CoT [Yao <i>et al.</i> , 2023]	67.1	-10.0%	68.5	-11.1%
SPP-Profile [Wang <i>et al.</i> , 2023c]	79.1	+5.9%	83.0	+7.8%
SPP [Wang <i>et al.</i> , 2023c]	79.9	+7.1%	84.7	+10.0%
<b>AutoAgents</b>	<b>82.0</b>	<b>+9.9%</b>	<b>85.9</b>	<b>+11.6%</b>



# 8.1 Further Analysis



Ablationsstudie: Komponentenwichtigkeit

Eine detaillierte Ablationsstudie an 20 Instanzen der Trivia Creative Writing-Aufgabe offenbart die Bedeutung jeder Komponente

Methods	N (# trivia questions) = 5	
	Score (%)	$\Delta$ (v.s Standard %)
Standard	74.6	0.0%
CoT [Yao <i>et al.</i> , 2023]	66.0	-11.5%
SPP-Profile [Wang <i>et al.</i> , 2023c]	74.0	-0.01%
SPP [Wang <i>et al.</i> , 2023c]	84.4	+13.1%
<b>AutoAgents</b> w/o observers	87.0	<b>+16.6%</b>
<b>AutoAgents</b> w/o self-refinement	87.0	<b>+16.6%</b>
<b>AutoAgents</b> w/o dynamic memory	89.0	<b>+19.3%</b>
<b>AutoAgents</b>	<b>90.0</b>	<b>+20.6%</b>



## 8.2 Vergleich mit bestehenden Frameworks



AutoAgents unterscheidet sich fundamental von bestehenden Multi-Agent-Frameworks durch mehrere Schlüsselinnovationen

Framework	Dynamic Agent Generation Method	Number of Agent	Multi-agent Conversation	Self-Refinement Action	Collaborative Refinement Action
AutoGPT [Gravitas, 2023]	X	1	X	✓	X
BabyAGI [Nakajima, 2023]	X	3	✓	X	X
Generative Agents [Park <i>et al.</i> , 2023]	X	25	✓	✓	X
Camel [Li <i>et al.</i> , 2023]	X	2	✓	X	X
GPT-bargaining [Fu <i>et al.</i> , 2023]	X	3	✓	✓	X
MetaGPT [Hong <i>et al.</i> , 2023]	X	Unlimited	✓	X	X
AutoGen [Wu <i>et al.</i> , 2023]	X	Unlimited	✓	X	X
Social Simulacra [Park <i>et al.</i> , 2022]	Single Agent	Unlimited	✓	X	X
Epidemic Modeling [Williams <i>et al.</i> , 2023]	Single Agent	Unlimited	✓	X	X
ExpertPrompting [Xu <i>et al.</i> , 2023]	Single Agent	1	X	X	X
SSP [Wang <i>et al.</i> , 2023c]	Single Agent	Unlimited	✓	X	X
AgentVerse [Chen <i>et al.</i> , 2023a]	Single Agent	Unlimited	✓	X	X
<b>AutoAgents</b>	Multi-agent Discussion	Unlimited	✓	✓	✓



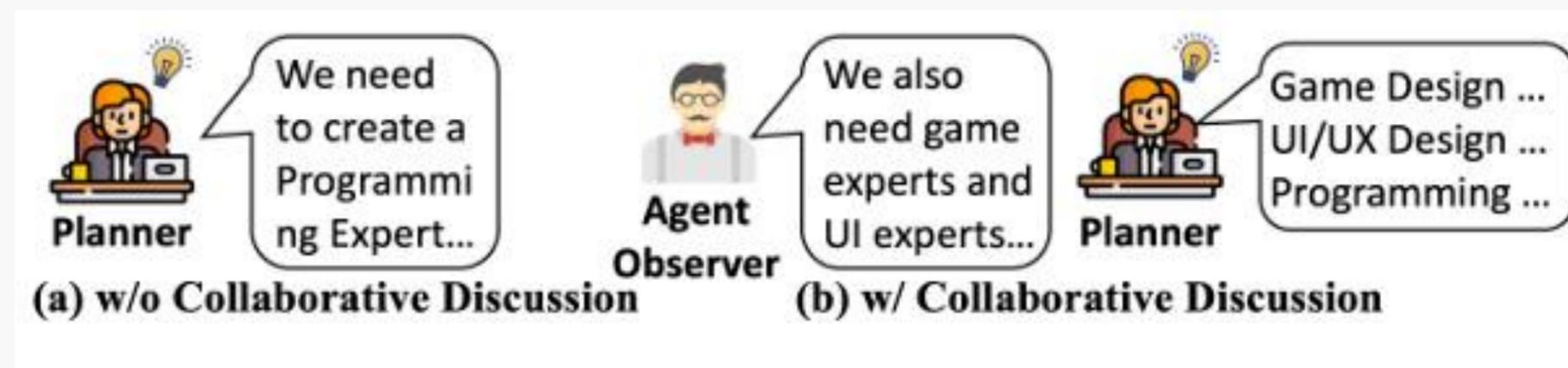
# 9. Praktische Anwendungen und Fallstudien



## Software-Entwicklung

Eine Fallstudie zur Entwicklung eines Python-basierten Tetris-Spiels zeigt die Überlegenheit kollaborativer Diskussion in der Entwurfsphase

- Ohne Observer: Der Planner generiert ausschließlich Programmierer und vernachlässigt den ganzheitlichen Prozess
- Mit Observer: Einbeziehung von Spieldesign-Experten, UI-Design-Experten und Test-Experten



# Conclusion



AutoAgents präsentiert ein innovatives Framework für dynamische Multi-Agenten-Systeme auf Basis großer Sprachmodelle. Es kombiniert flexible Agentengenerierung, kollaborative Planung und effektive Beobachtermechanismen, um komplexe Aufgaben effizient zu lösen. Die Ergebnisse zeigen klare Vorteile gegenüber bisherigen Ansätzen und machen AutoAgents zu einem wichtigen Schritt in Richtung leistungsfähiger, adaptiver KI-Systeme.





*Danke für Ihre Aufmerksamkeit*

