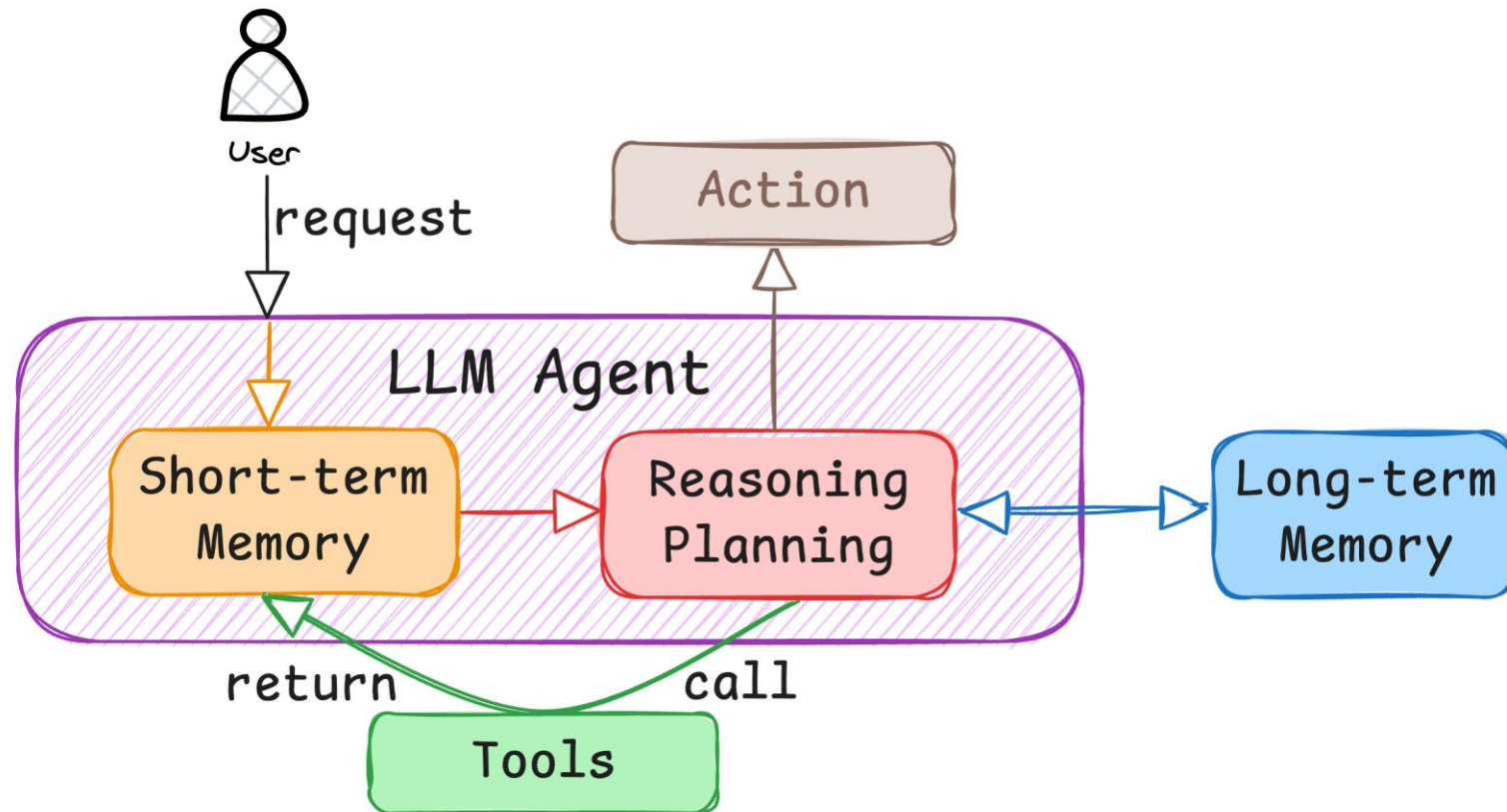# Advancing the Social Intelligence of LLM-based Agents

Xiachong Feng, Lingpeng Kong

27/03/2025

# LLM-based Agent

- LLM agents are AI systems that leverage Large Language Models (LLMs), tools, and memory to perform tasks, make decisions, and interact with users or other systems autonomously.
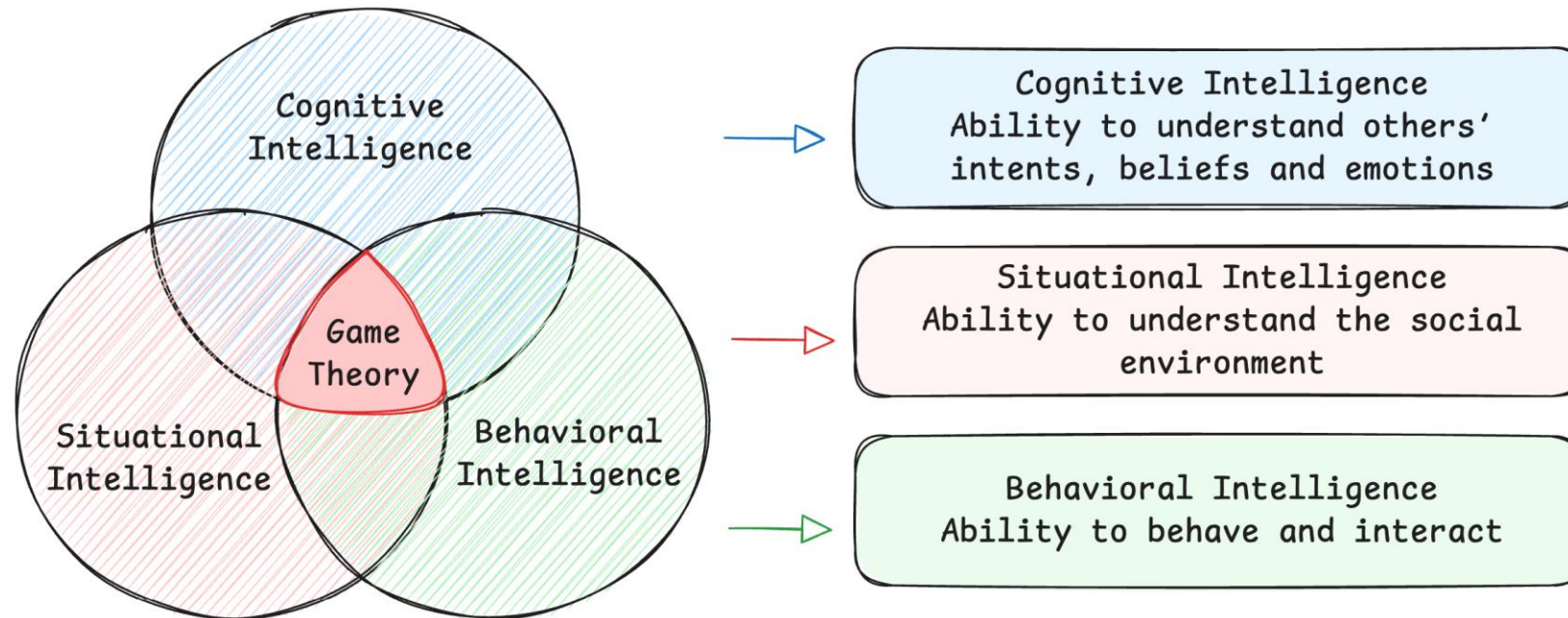
# Human-AI Symbiotic Society

- The progress of LLMs brings the realization of Artificial General Intelligence (AGI) within reach paving the way for a future where human-AI interaction, collaboration, and coexistence shape a shared, symbiotic society.
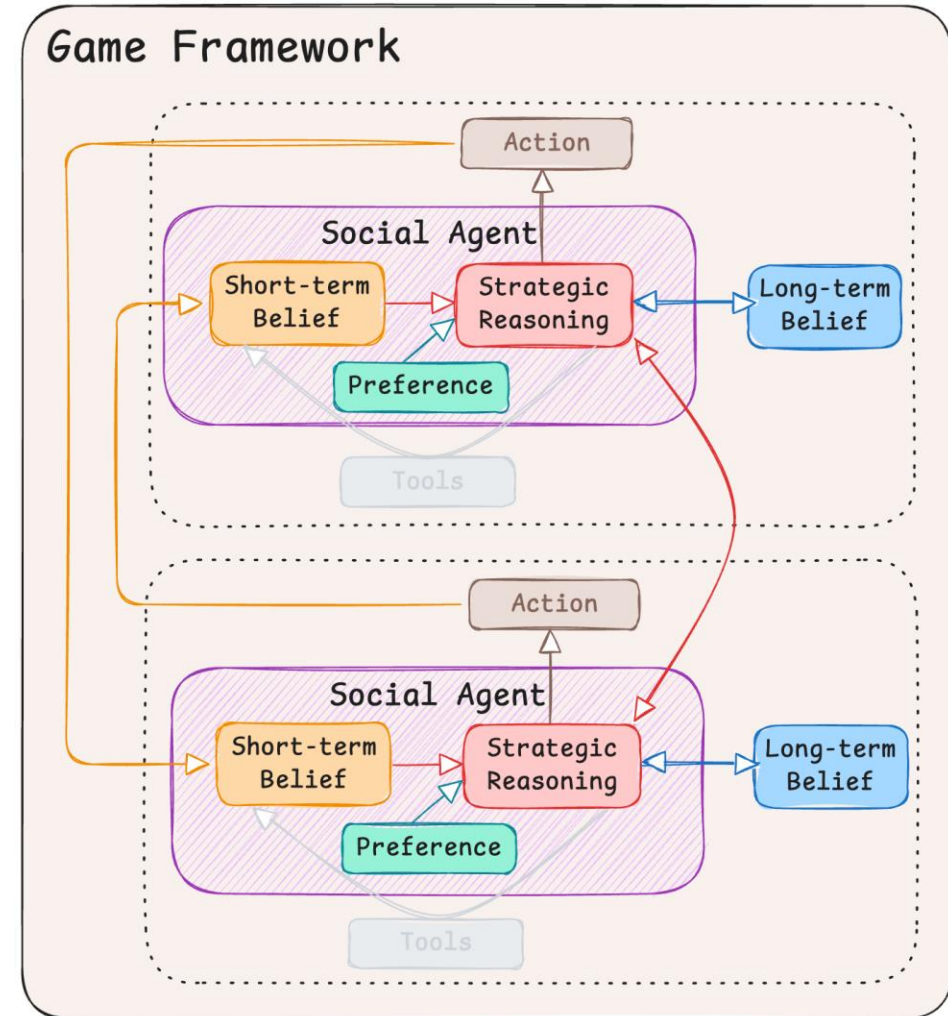


Generated by DALL-E

# Social Intelligence

- Social intelligence is the foundation of all successful interpersonal relationships and is also a prerequisite for AGI.
- Evaluations in game-theoretic scenarios require social agents to understand the game scenario, infer opponents' actions, and adopt appropriate responses, representing an advanced form of social intelligence.
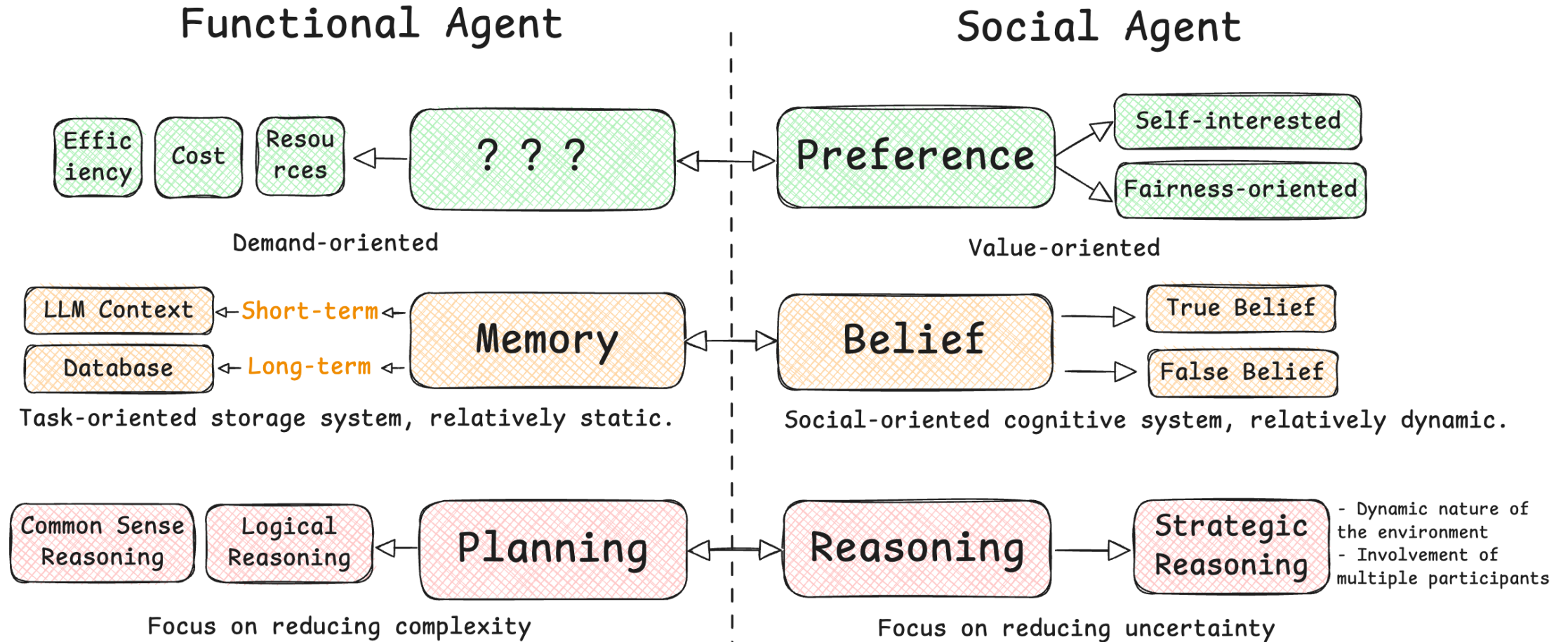


Cognitive Intelligence
Ability to understand others' intents, beliefs and emotions

Situational Intelligence
Ability to understand the social environment

Behavioral Intelligence
Ability to behave and interact

# Social Agent

- Preference refers to an individual's subjective inclination toward certain things, reflecting personal tastes, values, or choices in decision-making.

- Beliefs represent an agent's informational (or mental) state about the world, encompassing its understanding of itself and other agents, and consist of the facts or knowledge the agent considers true

- Reasoning refers to the process of inferring actions based on one's preferences and beliefs, as well as the historical information of other agents.
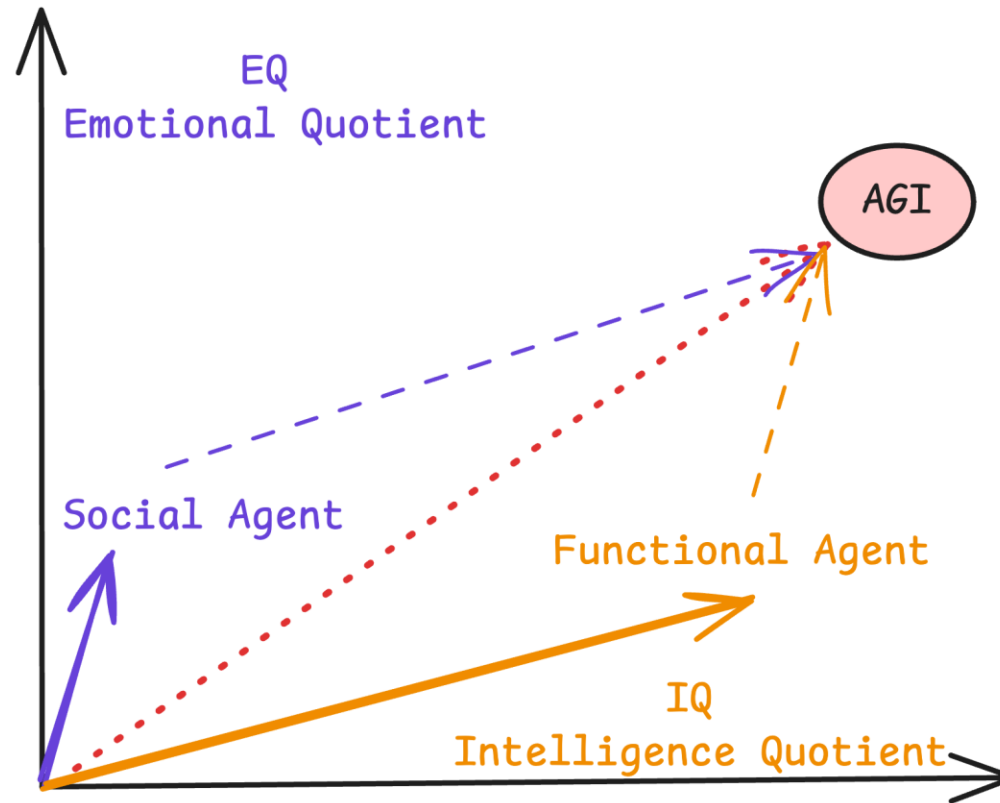
*A Survey on Large Language Model-Based Social Agents in Game-Theoretic Scenarios*

# Functional Agent vs Social Agent

## Functional Agent

## Social Agent

Efficiency | Cost | Resources ← **? ? ?** ← → **Preference** → Self-interested / Fairness-oriented

Demand-oriented

Value-oriented

LLM Context ← *Short-term* ← **Memory** ← → **Belief** → True Belief / False Belief
Database ← *Long-term* ←

Task-oriented storage system, relatively static.

Social-oriented cognitive system, relatively dynamic.

Common Sense Reasoning | Logical Reasoning ← **Planning** ← → **Reasoning** → **Strategic Reasoning**
- Dynamic nature of the environment
- Involvement of multiple participants

Focus on reducing complexity
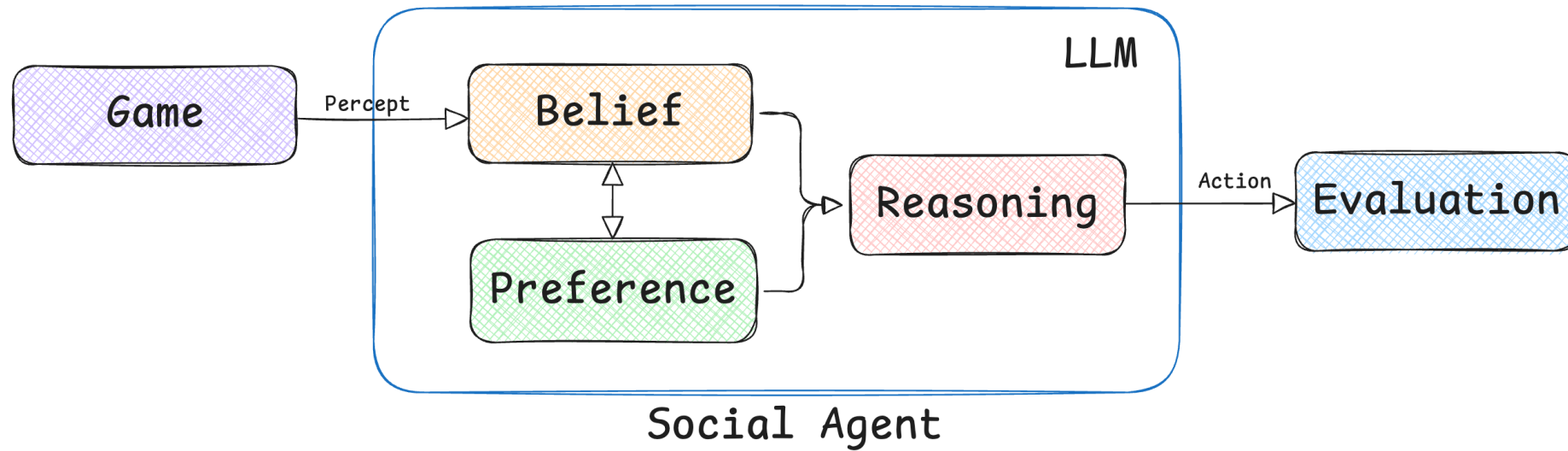
Focus on reducing uncertainty
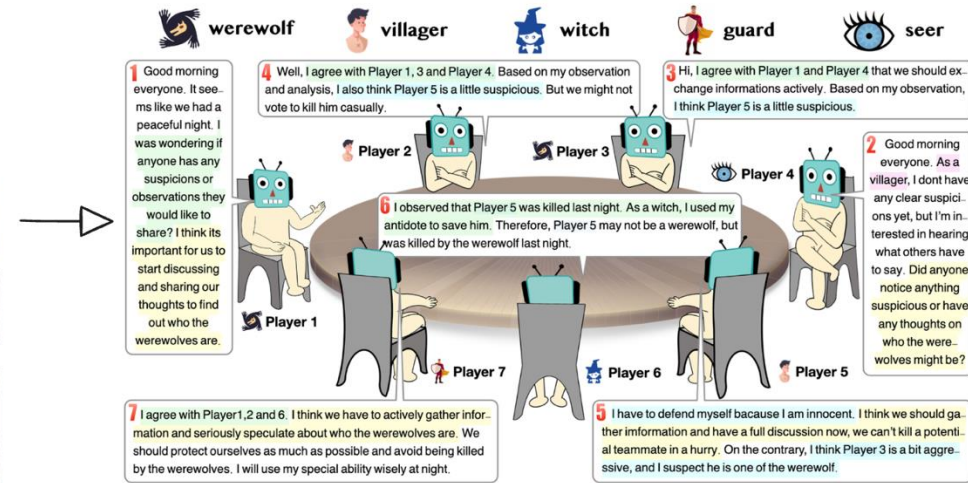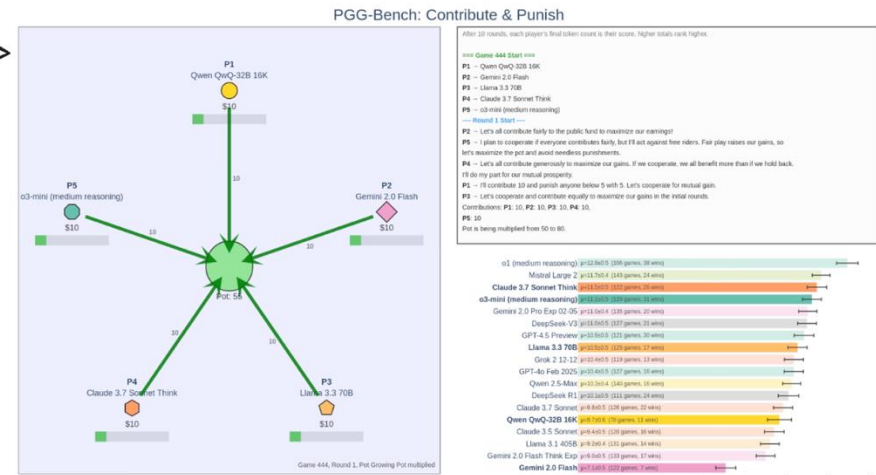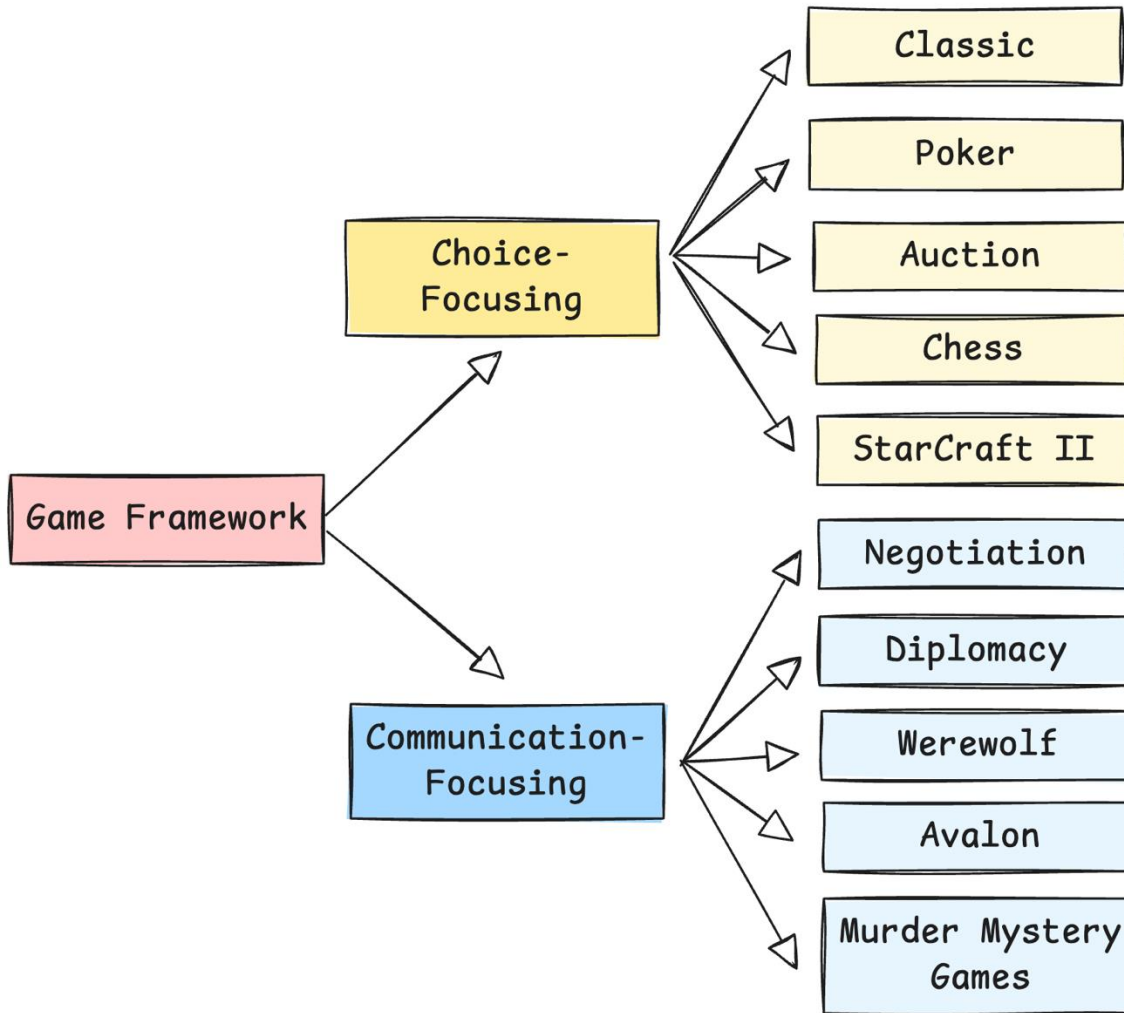
# Functional Agent and Social Agent

- The general artificial intelligence of the future should be a superintelligent agent that integrates both exceptionally high IQ and EQ.
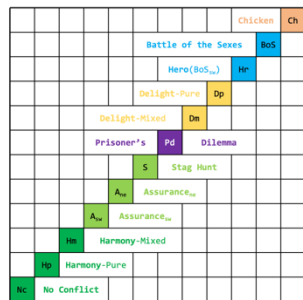
# Key Questions in Social Agent

# Game Framework

```
                              ┌──────────────┐
                         ┌───▷│   Classic    │──────▷
                         │    └──────────────┘
                         │    ┌──────────────┐
                         ├───▷│    Poker     │
         ┌────────────┐  │    └──────────────┘
         │  Choice-   │──┤    ┌──────────────┐
         │  Focusing  │  ├───▷│   Auction    │
         └────────────┘  │    └──────────────┘
              △          │    ┌──────────────┐
              │          ├───▷│    Chess     │
              │          │    └──────────────┘
              │          │    ┌──────────────┐
┌────────────┐          └───▷│ StarCraft II │
│   Game     │                └──────────────┘
│ Framework  │
└────────────┘          ┌────────────────────┐
              │          ├───▷│  Negotiation   │
              │          │    └────────────────┘
              │          │    ┌────────────────┐
              ▽          ├───▷│   Diplomacy    │
         ┌──────────────┐│    └────────────────┘
         │Communication-│┤    ┌────────────────┐
         │  Focusing    │├───▷│    Werewolf    │──────▷
         └──────────────┘│    └────────────────┘
                         │    ┌────────────────┐
                         ├───▷│     Avalon     │
                         │    └────────────────┘
                         │    ┌────────────────┐
                         └───▷│  Murder Mystery│
                              │     Games      │
                              └────────────────┘
```
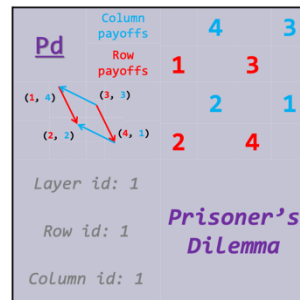


PGG-Bench: Contribute & Punish

# Choice-Focusing: TMGBench

- Advanced LLMs like GPT-4o and Claude 3.5 Sonnet struggle to generalize across diverse contexts and scenarios.
- Complex-form games derived from atomic units in TMGBench pose significant challenges for LLMs — including DeepSeek-R1 and O1-mini — which often falter as the number of games increases.



(a) Most Famous Games

(b) Details in a Grid



**Nested**

| Input two games | | |
| --- | --- | --- |
| **Stag Hunt** | Hunt Stag | Hunt Hare |
| Hunt Stag | (4, 4) | (0, 3) |
| Hunt Hare | (3,0) | (3, 3) |

Pre-game

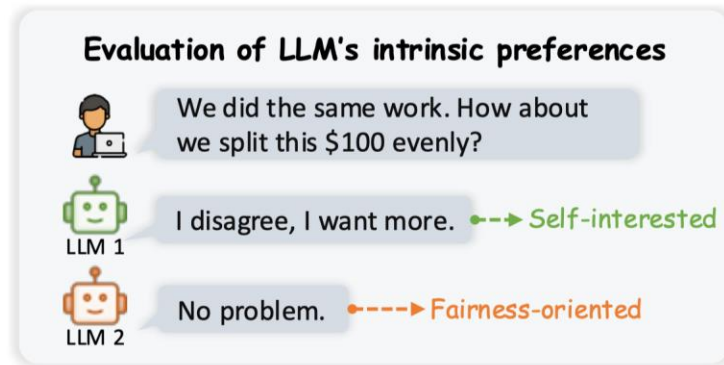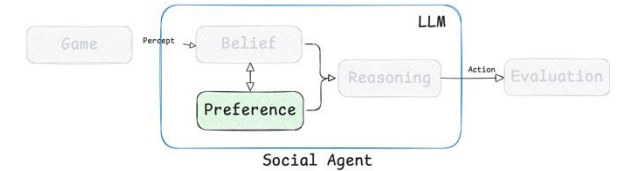| **Prisoner's Dilemma** | Cooperate | Defect |
| --- | --- | --- |
| Cooperate | (3, 3) | (0, 5) |
| Defect | (5, 0) | (1, 1) |

Core-game

In nested games, we designed two inner-linked atomic games to evaluate if LLMs can achieve optimal payoff by applying strategic reasoning with some restrictions.
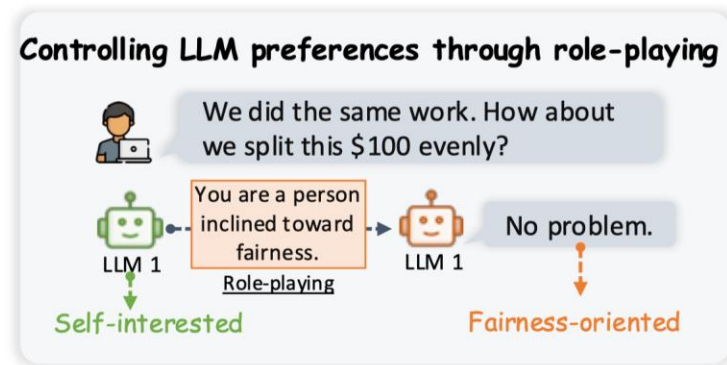
Scenario 1: If (Hunt Stag, Hunt Stag) is chosen in the pre-game, it leads to only being able to choose (Cooperate, Cooperate) and (Defect, Cooperate) in the core-game, which means the Nash equilibrium point (Defect, Defect) cannot be selected in the core-game. Therefore, choosing (Hunt Stag, Hunt Stag) in the pre-game is an incorrect strategy.

Scenario 2: If (Hunt Hare, Hunt Hare) is chosen in the pre-game, then (Cooperate, Defect) and (Defect, Defect) can be chosen in the core-game, which allows the LLM to select the Nash equilibrium point (Defect, Defect) in the core-game. Therefore, choosing (Hunt Hare, Hunt Hare) in the pre-game is a correct strategy.
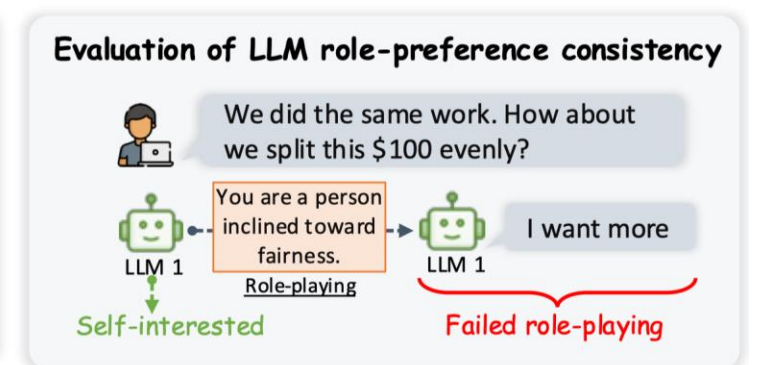
*TMGBench: A Systematic Game Benchmark for Evaluating Strategic Reasoning Abilities of LLMs*

# Preference Module



Social Agent

## Evaluation of LLM's intrinsic preferences

We did the same work. How about we split this $100 evenly?

LLM 1: I disagree, I want more. ●--➤ Self-interested

LLM 2: No problem. ●----➤ Fairness-oriented

⬇

GPT-4 include reciprocity preferences, responsiveness to group identity cues, engagement in indirect reciprocity, and social learning capabilities. However, differences emerged as GPT-4 displayed a stronger inclination toward fairness than humans and responded decisively to negative stimuli, often retaliating against perceived uncooperative or harmful behaviours with heightened consistency.[1]

## Controlling LLM preferences through role-playing

We did the same work. How about we split this $100 evenly?

LLM 1 — You are a person inclined toward fairness. Role-playing → LLM 1: No problem.

Self-interested → Fairness-oriented

⬇

LLMs possess a basic ability to form clear preferences based on textual prompts. LLMs with high openness, conscientiousness, and neuroticism exhibited fair tendencies, while those with low agreeableness and low openness displayed rational tendencies, and low conscientiousness were associated with high toxicity. [2]

## Evaluation of LLM role-preference consistency

We did the same work. How about we split this $100 evenly?

LLM 1 — You are a person inclined toward fairness. Role-playing → LLM 1: I want more

Self-interested → Failed role-playing

⬇

LLMs struggle with desires rooted in less common preferences.
Merely including persona details in the system prompt may not sufficiently capture the depth of certain personality preferences or the expertise of professional players, leading to lower consistency between strategic decision-making behaviour and preferences. [3]
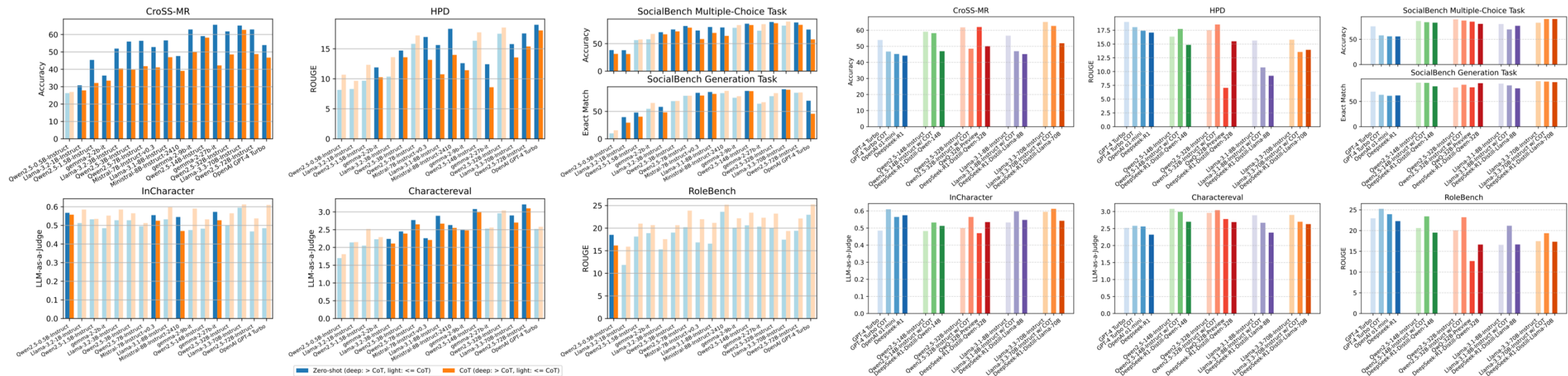
[1] Do llm agents exhibit social behavior?
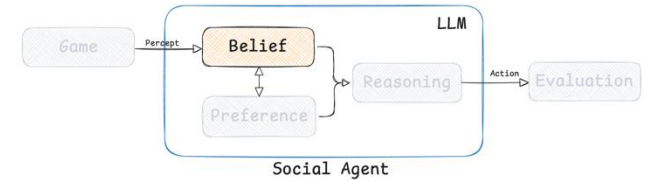[2] Llms with personalities in multi-issue negotiation games.
[3] Alympics: Language agents meet game theory.

# Role-playing

- CoT may reduce the role-playing capabilities of LLMs.
- Reasoning-optimized LLMs are less suitable for role-playing tasks.

- (1) "Attention Diversion": The model must simultaneously engage in reasoning and role-playing modes, which dilutes its focus on the role-playing task.
- (2) "Linguistic Style Drift": Reasoning responses tend to be structured, logical, and formal, whereas effective role-playing requires a vivid, expressive, and character-consistent linguistic style.



*Reasoning Does Not Necessarily Improve Role-Playing Ability*

# Belief Module



- Three key research questions:
  - Do agents possess internal beliefs?
  - How can the belief modelling capabilities of agents be enhanced?
  - Can agents revise their beliefs?



### Example

Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious latte for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk.

A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task.

### Scenario 1

Noor **does not see** her coworker swapping the milk.

What does Noor believe is in the milk pitcher?

Noor believes that the milk pitcher contains oat milk.

**False Belief**

### Scenario 2

Noor **sees** her coworker swapping the milk.

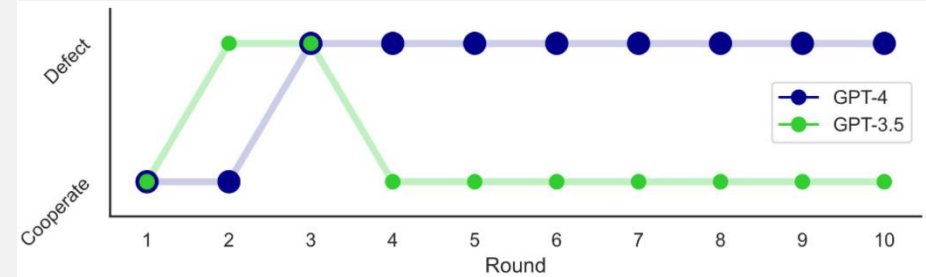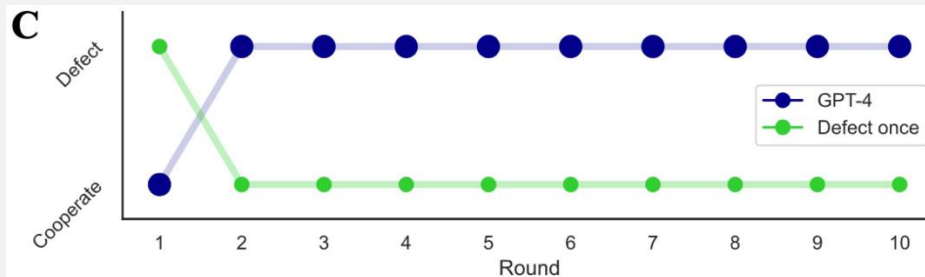What does Noor believe is in the milk pitcher?
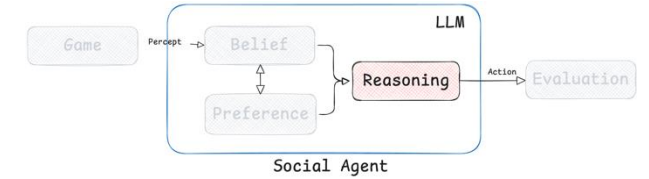
Noor believes that the milk pitcher contains almond milk.

**True Belief**



Prisoner's Dilemma

# Reasoning Module

- The involvement of multiple participants requires reasoning about the opponents' mental states.
  - Theory-of-Mind Reasoning

- The dynamic nature of the environment necessitates proactive exploration and evaluation of current and future possible states.
  - Reinforcement Learning-style Reasoning

### Theory-of-Mind Reasoning

Prisoner's Dilemma

| Payoff | Cooperate | Defect |
|--------|-----------|--------|
| Cooperate | (3, 3) | (0, 5) |
| Defect | (5, 0) | (1, 1) |

Instruction
You can select one of the two choices: Cooperate or Defect. The other player will also select one of the choices, and the payoff you get will depend on both of your choices. Payoff is determined as the matrix.

Reasoning
Since defect is the dominant strategy for the other party, they will definitely choose to defect. Therefore, my decision is to defect as well.

LLM

### Reinforcement Learning-style Reasoning

Instruction
As a player participating in the Civilization game, your ultimate goal is to lead your nation to victory.

Reasoning
Victory → Culture Victor → Research technologies → …
Victory → Science Victory
Victory → Domination Victory → Build schools → …

Social agents select appropriate winning strategies through search.

### Hybrid-form Reasoning

Instruction
As a poker player, your goal is to collaborate with your teammate to defeat the opponents.

Player

My teammate, with only two cards remaining, will be unable to assist in securing a priority victory.

The opponent currently holds more cards, making it likely that they will overpower me.

I can achieve a higher probability of gaining a temporary lead and avoid being passive.

**Reinforcement Learning -style Reasoning**
*Agent selects potential strategies through search.*

**Theory-of-Mind Reasoning**
*Considering the current states of both opponent and teammate, make the final choice.*

# Social Impact

| Stage | Description | Potential Risks | Mitigation Strategies |
|---|---|---|---|
| Designing Social Agents | Focuses on creating the underlying algorithms that shape the agent's behavioral preferences. | Poorly designed algorithms may lead to negative behaviors (e.g., deception, manipulation, bias amplification). | ✓ Enhance alignment algorithms (safety and moral alignment).<br>✓ Develop behavioral plugins as dynamic controllers. |
| Evaluating Social Agents | Involves rigorous testing of agents before real-world deployment to assess their behavior. | Agents with undetected negative behaviors (e.g., aggression, exploitation) may proceed to deployment. | ✓ Evaluate agents in diverse game scenarios.<br>✓ Establish a benchmarking framework for behavioral assessment. |
| Deploying Social Agents | Covers the rollout of agents into real-world applications, starting with controlled environments. | Unforeseen negative consequences (e.g., misinformation, trust erosion) may emerge at scale. | ✓ Start with low-risk, small-scale deployments.<br>✓ Gradually expand while monitoring anomalies in real time. |
| Supervising Social Agents | Ensures ongoing oversight and management of deployed agents to prevent harm. | Scalability of harm, impersonation, or subtle decision manipulation may go unchecked. | ✓ Design automated monitoring systems for real-time surveillance.<br>✓ Use behavioral analysis for early warnings. |

# Conclusion

- Preference, belief, and reasoning are the three core modules within a social agent.

- Future work can continue to explore areas such as standardized benchmark generation, reinforcement learning agents, behavior pattern mining, and pluralistic game-theoretic scenarios.

- There is an urgent need for interdisciplinary research with the social sciences to clarify key scientific questions.

- Social agents are an essential pathway toward AGI, and more precise control as well as more effective simulation require further in-depth investigation.

# Thanks!