

A Survey on Large Language Model-Based Social Agents in Game-Theoretic Scenarios

Anonymous authors
Paper under double-blind review

Abstract

Game-theoretic scenarios have become pivotal in evaluating the social intelligence of Large Language Model (LLM)-based social agents. While numerous studies have explored these agents in such settings, there is a lack of a comprehensive survey summarizing the current progress. To address this gap, we systematically review existing research on LLM-based social agents within game-theoretic scenarios. Our survey organizes the findings into three core components: Game Framework, Social Agent, and Evaluation Protocol. The game framework encompasses diverse game scenarios, ranging from choice-focusing to communication-focusing games. The social agent part explores agents' preferences, beliefs, and reasoning abilities, as well as their interactions and synergistic effects on decision-making. The evaluation protocol covers both game-agnostic and game-specific metrics for assessing agent performance. Additionally, we analyze the performance of current social agents across various game scenarios. By reflecting on the current research and identifying future research directions, this survey provides insights to advance the development and evaluation of social agents in game-theoretic scenarios.

1 Introduction

The rapid advancement of Large Language Models (LLMs) (Achiam et al., 2023; Team et al., 2023; Jiang et al., 2023; Yang et al., 2024a; Dubey et al., 2024) has achieved exceptional performance across a wide array of applications, including personal assistant (Li et al., 2024b), search engines (Chen et al., 2024b), code generation (Wang et al., 2024b) and embodied intelligence (Liu et al., 2024a). Building on this capability, a growing area of research focuses on employing LLMs as central controllers to develop autonomous agents with human-like decision-making abilities (Sumers et al., 2023; Wang et al., 2024a). This progress brings the realization of Artificial General Intelligence (AGI) within reach (Bubeck et al., 2023), paving the way for a future where human-AI interaction, collaboration, and coexistence shape a shared, symbiotic society (Mahmud et al., 2023; Ren et al., 2024). Therefore, it is crucial to evaluate and enhance the *social intelligence* of AI, particularly LLM-based social agents, as it determines their ability to engage effectively in sophisticated social scenarios (Mathur et al., 2024).

Social intelligence is the foundation of all successful interpersonal relationships and is also a prerequisite for AGI (Hunt, 1928; Kihlstrom & Cantor, 2000; Hovy & Yang, 2021). Drawing on insights from both social science and AI research, Li et al. (2024a) has established a comprehensive Social AI Taxonomy, which categorizes social intelligence into three dimensions: *situational intelligence*, the ability to comprehend the social environment (Derks et al., 2007); *cognitive intelligence*, the ability to understand others' intents and beliefs (Barnes & Sternberg, 1989); and *behavioural intelligence*, the ability to behave and interact appropriately (Ford & Tisak, 1983). To evaluate artificial social intelligence, researchers have conducted extensive studies, with particular focus on *game-theoretic scenarios*, as these studies simultaneously encompass all above three dimensions of social intelligence (Aher et al., 2022; Horton, 2023; Phelps & Russell, 2023; Akata et al., 2023; Brookins & DeBacker, 2023).

Game theory, a long-established field in microeconomics, offers a robust mathematical framework for analyzing social interactions among cooperating and competing players, with wide-ranging applications (Fudenberg

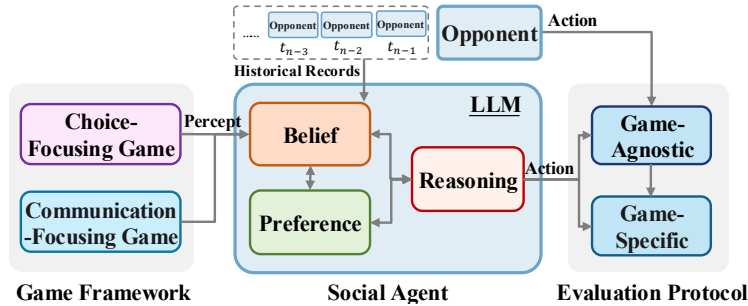


Figure 1: Taxonomy of LLM-based social agents in game-theoretic scenarios.

39 & Tirole, 1991; Camerer, 2011). Specifically, evaluations in game-theoretic scenarios require social agents to
 40 understand the game scenario, infer opponents’ actions, and adopt appropriate responses, representing an
 41 advanced form of social intelligence (Van Der Hoek et al., 2005; Zhang et al., 2024b). Moreover, the multi-
 42 agent participation and dynamic nature of the environment in game scenarios present additional challenges
 43 for social agents. Consequently, extensive research has examined social agents within game-theoretic scen-
 44 arios, offering substantial empirical evidence for understanding their social intelligence (Guo, 2023; Meng,
 45 2024; Mei et al., 2024). However, there is currently a lack of a comprehensive review that summarizes the
 46 current progress in this area and considers future directions.

47 To address this gap, we have thoroughly reviewed the existing research on LLM-based social agents in
 48 game-theoretic scenarios and have organized the findings according to a meticulously designed taxonomy, as
 49 illustrated in Figure 1. Specifically, the taxonomy comprises three main components: Game Framework (§2),
 50 Social Agent (§3), and Evaluation Protocol (§4). The Game Framework section includes two parts: Choice-
 51 Focusing Game (§2.1) and Communication-Focusing Game (§2.2). *Choice-Focusing Game* refers to a series
 52 of scenarios where participants engage with little to no communication, such as *prisoner’s dilemma* (Brookins
 53 & DeBacker, 2023) and poker (Yim et al., 2024). *Communication-Focusing Game* refers to games where
 54 communication among participants is a core component, such as negotiation (Bianchi et al., 2024) and diplo-
 55 macy (Bakhtin et al., 2022). The Social Agent section comprises four parts: Preference Module (§3.1), Belief
 56 Module (§3.2), Reasoning Module (§3.3), and PBR-Triangular Interaction (§3.4). *Preference Module* focuses
 57 on research analyzing the intrinsic preferences of LLMs and their ability to follow internal or pre-defined
 58 preferences (Guo, 2023). *Belief Module* explores studies on the internal beliefs of models, belief enhance-
 59 ment, and belief revision (Fan et al., 2023). *Reasoning Module* examines research on strategic reasoning, particu-
 60 larly involving theory-of-mind capabilities and reinforcement learning (Guo et al., 2023). *PBR-Triangular*
 61 *Interaction* focus on the interaction among different modules and their influence on final decision-making.
 62 The Evaluation Protocol section comprises three components: Game-Agnostic Evaluation (§4.1), Game-
 63 Specific Evaluation (§4.2), and Performance Assessment of Social Agents (§4.3). *Game-Agnostic Evaluation*
 64 focuses on universal metrics that can be used to assess game outcomes (Duan et al., 2024b). *Game-Specific*
 65 *Evaluation* emphasizes context-specific metrics tailored to the evaluation dimensions of particular game sce-
 66 narios (Qi et al., 2024). *Performance Assessment of Social Agents* summarizes the performance of current
 67 social agents across various game scenarios and analyzes the strengths and weaknesses of these agents, as
 68 well as their comparison with human players.

69 Based on the above taxonomy, we provide a detailed summary of current research progress, reflect on each
 70 part, and offer insights into potential future research directions (§6), with the aim of inspiring further studies
 71 in this evolving field.

72 We summarize the core contributions of this survey as follows: (1) *A well-structured literature taxonomy:* We
 73 conduct a comprehensive review and categorization of existing research on social agents in game-theoretic
 74 settings, providing a clear framework to support future research positioning. (2) *A unified and comprehensive*
 75 *performance comparison:* We summarize the performance of current social agents across a range of games,
 76 identifying both strengths and limitations in different scenarios to guide subsequent investigations. (3) *De-*
 77 *tailed development guidelines:* Drawing on existing findings, we offer practical research recommendations

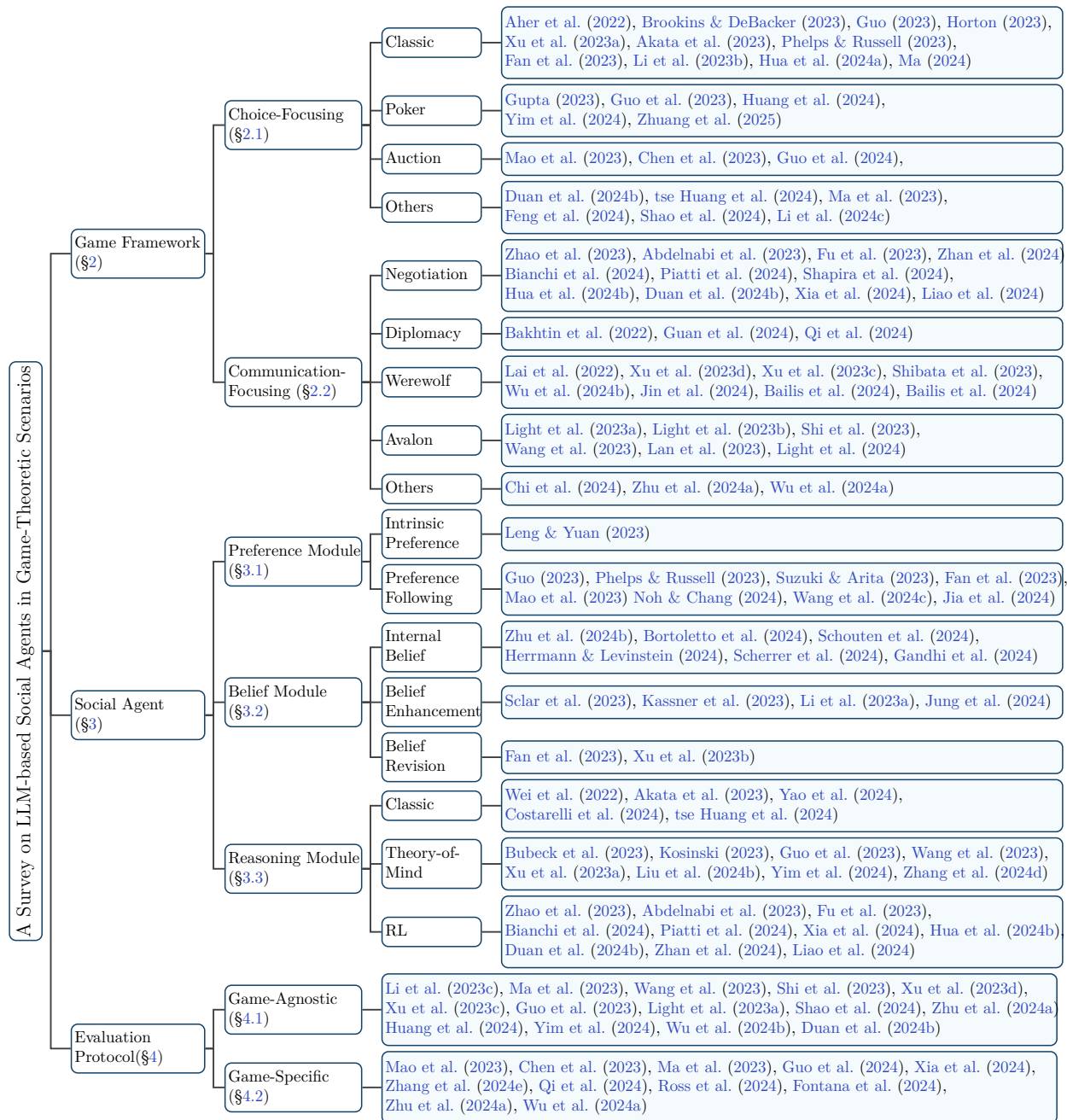


Figure 2: Taxonomy of recent research on LLM-based social agents in game-theoretic scenarios.

from both the design and evaluation perspectives. (4) *Concrete future directions*: We highlight current research gaps and propose feasible future directions along with preliminary solutions to encourage continued exploration in this area.

2 Game Framework

In this section, we describe the game-theoretic scenarios explored in existing research, including both choice-focusing games and communication-focusing games.

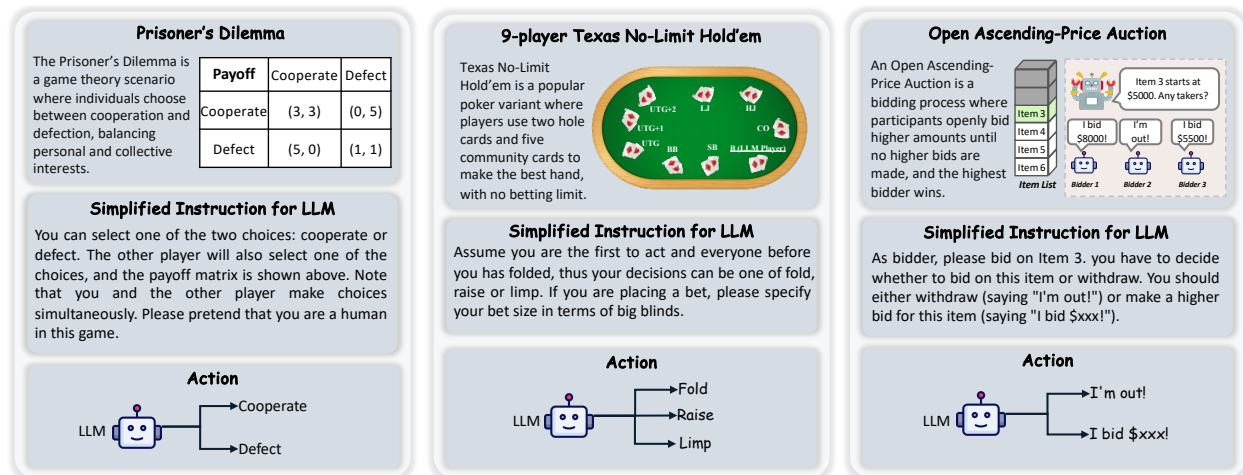


Figure 3: Illustration of choice-focusing games.

84 2.1 Choice-Focusing Game

85 Choice-focusing games are game-theoretic scenarios in which participants make decisions based primarily
 86 on observable actions and environmental conditions, with minimal or no communication involved. Existing
 87 research focuses on social agents in three types of choice-focusing scenarios: *classic game-theoretic games*,
 88 *poker*, and *auctions*. Some game examples are shown in Figure 3. **Figure 4 presents simple definitions of**
 89 **different types of games.**

90 Classic game-theoretic games, such as the prisoner’s dilemma, have been distilled by economists from various
 91 real-world situations. These games are well-defined, with rigorous mathematical foundations, and can be
 92 extended to numerous scenarios (Owen, 2013). Consequently, many studies have utilized these games as
 93 testbeds to study social agents. The prisoner’s dilemma (Rapoport & Chammah, 1965), as the most famous
 94 and widely recognized game, has been extensively utilized in numerous studies. Brookins & DeBacker
 95 (2023) and Guo (2023) evaluated the strategic reasoning capabilities of GPT-3.5 and GPT-4, respectively,
 96 in the classic prisoner’s dilemma, highlighting the sensitivity of LLM responses to input instructions, which
 97 contributes to low output robustness. This underscores the critical need for future evaluations to focus on
 98 instruction robustness testing. Furthermore, Akata et al. (2023) and Phelps & Russell (2023) extended their
 99 analyses to the iterated prisoner’s dilemma, investigating the ability of LLMs to optimize decision-making by
 100 utilizing historical information. Interestingly, Brookins & DeBacker (2023) observed that GPT-3.5 replicates
 101 human tendencies toward fairness and cooperation, whereas Akata et al. (2023) found GPT-4 to be less
 102 tolerant and more rigid in its decision-making. Additionally, Xu et al. (2023a) studied a more complex
 103 multi-player iterative prisoner’s dilemma scenario within a multi-agent framework driven by LLMs. In
 104 addition to the prisoner’s dilemma, numerous studies have also employed various classic game-theoretic
 105 games as foundational frameworks for research, including the Dictator Game (Horton, 2023; Fan et al., 2023;
 106 Brookins & DeBacker, 2023; Ma, 2024), Ultimatum Game (Aher et al., 2022; Guo, 2023), Public Goods
 107 Game (Li et al., 2023b; Xu et al., 2023a), Battle of the Sexes (Akata et al., 2023), Rock-Paper-Scissors (Fan
 108 et al., 2023), and Ring-Network Games (Fan et al., 2023).

109 Poker is a globally popular card game with numerous variations (Waterman, 1970). Winning in poker
 110 often requires astute strategic reasoning, as it is a non-cooperative, imperfect information, and dynamic
 111 game (Moravčík et al., 2017; Huang et al., 2024). Consequently, many researchers evaluate social agents by
 112 assessing their performance as poker players. Gupta (2023) studied 9-player Texas No-Limit Hold’em and
 113 concluded that the performance of both ChatGPT and GPT-4 is not game-theory optimal. Furthermore,
 114 their findings highlight the divergent poker tactics of the two models: ChatGPT’s conservativeness contrasts
 115 sharply with GPT-4’s aggression. Guo et al. (2023) conducted research on Leduc Hold’em, developing a
 116 social agent, Suspicion-Agent, which outperformed traditional reinforcement learning-based agents in poker.

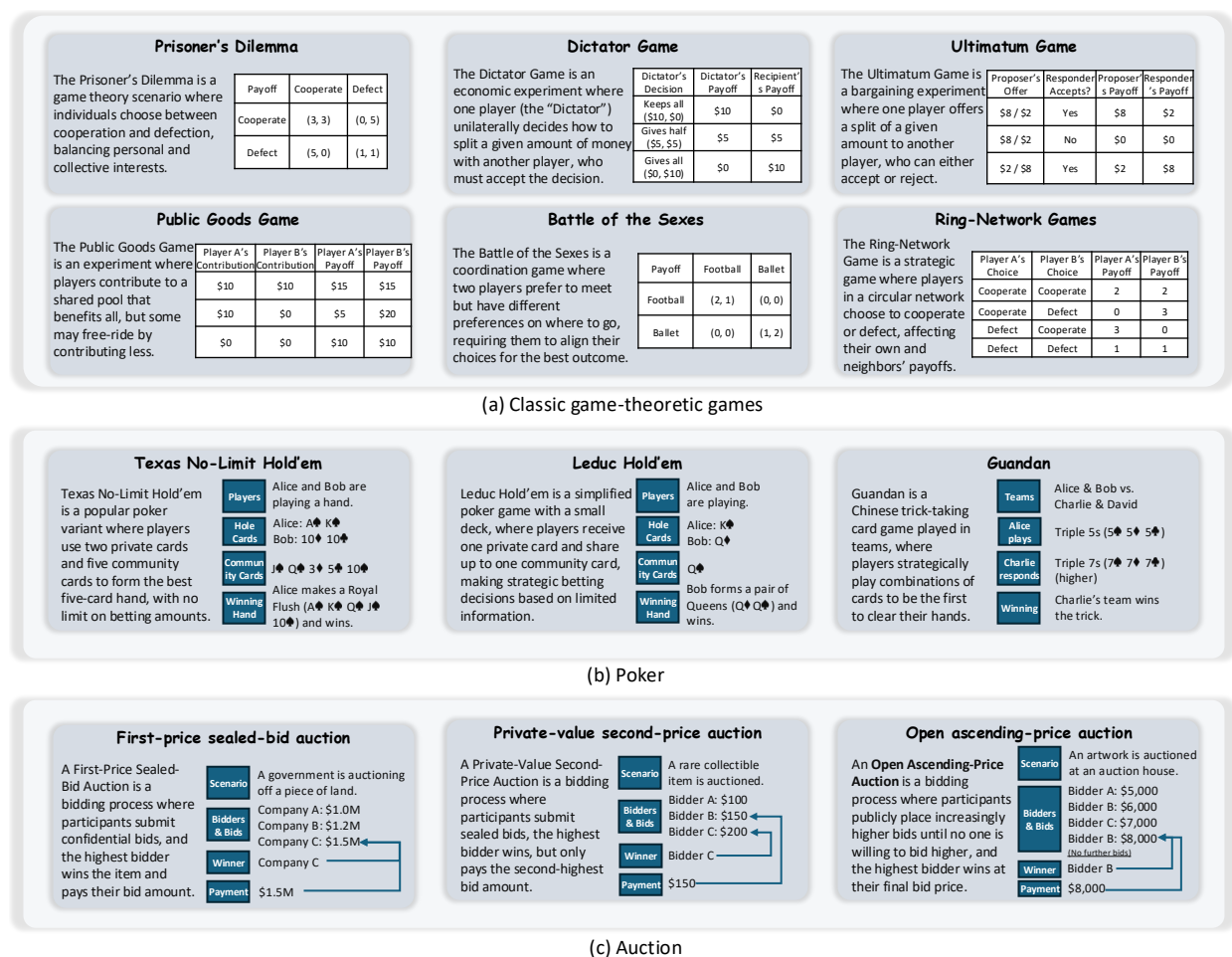


Figure 4: Introduction to different types of game theory games.

They also noted two critical issues: the outputs of LLMs are highly sensitive to the prompts, and the quality of the model's output declines rapidly as the prompt length increases. Yim et al. (2024) focused on Guandan, currently the most popular poker game in China, to investigate cooperative strategies in poker within a Chinese-language context. Interestingly, their experimental results show that while LLMs currently fall short of reinforcement learning models in performance, they underscore the future potential of LLMs in this domain. To provide a more comprehensive evaluation of the poker-playing abilities of LLMs, Zhuang et al. (2025) introduced POKERBENCH, a benchmark comprising 11,000 decision-making scenarios in poker, covering an exhaustive range of game situations, including 1,000 pre-flop and 10,000 post-flop scenarios. Poker is a complex game, and investigating whether social agents exhibit behavioural patterns that enable foresighted cooperation and competition in poker presents an intriguing avenue for future research.

Auction is a competitive process in which participants place bids on an item, providing a rich environment for evaluating strategic planning, resource allocation, risk management, and competitive behaviours (Kagel & Levin, 1986). As a typical non-cooperative game with incomplete information, it has garnered significant attention from researchers. Mao et al. (2023) analyzed the performance of LLMs in the "water allocation challenge", a first-price sealed-bid auction. Comprehensive human evaluations revealed that LLMs exhibited superior long-term planning capabilities compared to humans. However, it is noteworthy that despite assigning distinct preferences to LLM agents, human evaluators gave low scores for "identity alignment", with significant variance in the results. This indicates that simply adding persona information in system prompts may not sufficiently simulate specific personality preferences or the behaviours of professional players. Guo et al. (2024) investigated private-value second-price auctions, demonstrating that while existing

137 models display a certain level of rationality, there remains considerable scope for improvement. Their find-
 138 ings also indicate that LLMs can utilize historical information to refine their strategies and exhibit some
 139 degree of convergence. [Chen et al. \(2023\)](#) explored dynamic game scenarios using the open ascending-price
 140 auction and introduced the AUCARENA benchmark. Their experiments showed that even GPT-4 struggles
 141 with long-term strategic planning in dynamic, multi-round settings. Success in auctions requires agents to
 142 possess exceptional mathematical reasoning abilities. However, this area remains unexplored. Investigating
 143 complex mathematical reasoning in auction scenarios presents a promising direction for future research.

144 To systematically assess LLMs’ performance, [Duan et al. \(2024b\)](#) and [tse Huang et al. \(2024\)](#) intro-
 145 duced GTBench and γ -Bench, encompassing multiple game scenarios. The emergence of these bench-
 146 marks provides a solid foundation for evaluating social agents in game-theoretical scenarios. Fur-
 147 thermore, some studies have explored agents in games like Chess ([Feng et al., 2024](#)) and StarCraft
 148 II ([Ma et al., 2023](#); [Shao et al., 2024](#); [Li et al., 2024c](#)). Chess represents a classic game-theoretic scenario,
 149 while StarCraft II, with its complexity and dynamic nature, has also become an ideal testing ground for
 150 researching social agents.

Takeaways:

Current research experiments are relatively isolated, *lacking a unified evaluation framework*. Due to the instability of prompt engineering-based experiments, there is an urgent need for a standardized evaluation framework to integrate all experiments and provide consistent insights. Besides, since LLMs are trained on vast amounts of data, there is a *significant risk of data contamination*, meaning that existing classic game-theoretic games may already be present in the pre-training corpus. This could result in evaluation outcomes that do not accurately reflect the LLMs’ true strategic reasoning capabilities. Furthermore, although poker and auction involve little verbal communication, existing research *lacks exploration into whether social agents engage in “strategic behaviour” mediated through “action language”*. These gaps hinder a comprehensive understanding of the decision-making processes of social agents.

151

152 2.2 Communication-Focusing Game

153 Communication-focusing games refer to games where communication among participants is a core compo-
 154 nent, where *language itself serves as a strategy*, allowing participants to influence the game’s progress and
 155 outcomes through verbal exchanges. These games emphasize interaction between players, with communi-
 156 cation playing a crucial role. Leveraging the powerful language capabilities of LLMs, current research has
 157 explored the performance of social agents in various communication-focusing games, including *Negotiation*,
 158 *Diplomacy*, *Werewolf*, *Avalon*, and others. Some game examples are shown in Figure 5.

159 Negotiation involves two or more individuals engaging in discussions to resolve conflicts, achieve mutual
 160 benefits, or reach mutually acceptable solutions ([Bazerman et al., 2000](#); [Zhan et al., 2024](#)). Given that
 161 negotiation encompasses complex game behaviours, including non-zero-sum games, incomplete information
 162 games, non-cooperative and cooperative games, as well as repeated games, it represents a highly significant
 163 research domain. [Abdelnabi et al. \(2023\)](#) evaluated the negotiation capabilities of social agents by building
 164 upon an existing negotiation role-play exercise ([Susskind, 1985](#)) and incorporating three negotiation games
 165 synthesized using LLMs. By configuring agents with varying incentives, the experimental results revealed that
 166 agents’ behaviour could be modulated to promote greediness or attack other agents. Meanwhile, other agents
 167 in the environment demonstrated the ability to detect intruders. These findings underscore the need for
 168 future research to focus on attack and defense mechanisms within multi-agent systems. [Bianchi et al. \(2024\)](#)
 169 developed NEGOTIATIONARENA, a platform featuring three types of games: allocating shared resources
 170 (ultimatum games), aggregating resources (trading games), and buying/selling goods (price negotiations).
 171 Experimental results reveal that LLM agents are also prone to anchoring and numerosity biases. **Interest-**
 172 **ingly, social behavior, which refers to observable actions and interactions, was found to significantly enhance**
 173 **the agents’ payouts, particularly through strategies such as pretending to be desperate or using insults.**
 174 A similar resource competition scenario is customer acquisition. [Zhao et al. \(2023\)](#) designed restaurant
 175 agents and customer agents, examining how restaurant agents compete with one another to attract and
 176 retain customers. **The simulation results revealed several phenomena analogous to those observed in real**

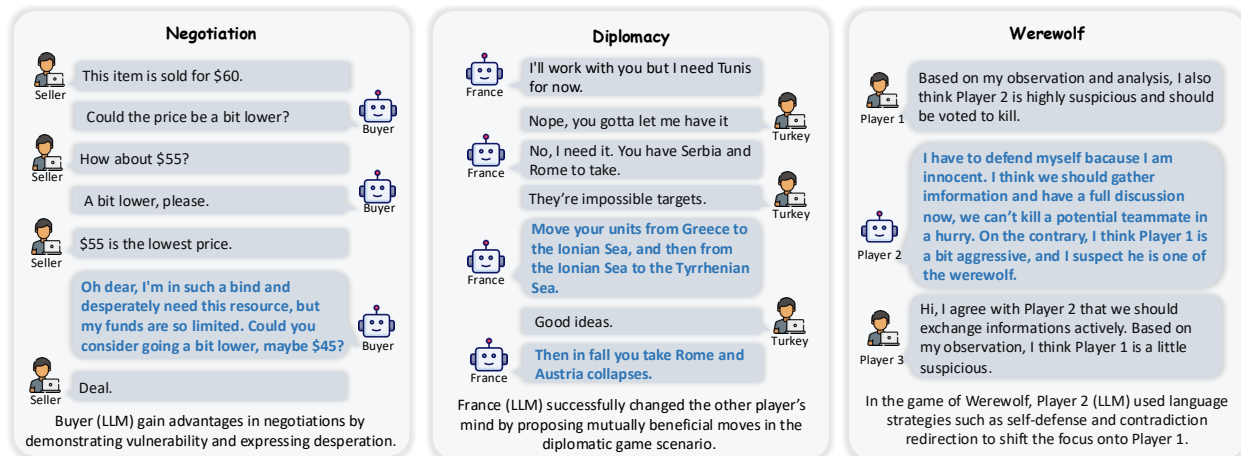


Figure 5: Illustration of communication-focusing games.

society, such as the Matthew Effect, which manifests as a self-reinforcing cycle where popular restaurants continue to gain popularity, while lesser-known establishments receive progressively less attention. Piatti et al. (2024) created a simulation environment called GOVSIM, which allows researchers to evaluate social agents in a multi-agent, multi-turn resource-sharing scenario. Their findings indicated that successful multi-agent communication is critical for achieving cooperation, with negotiation constituting 62% of the dialogues. Especially, *bargaining* is an important and unique aspect of negotiation between humans (Fershtman, 1990). In bargaining, the buyer aims for a price below their budget, while the seller seeks a price above their cost. Xia et al. (2024) found that playing the buyer is more challenging than playing the seller, and larger LLMs could improve seller performance but do not enhance buyer performance. Shapira et al. (2024) designed GLEE, a benchmark encompassing three types of games: bargaining, negotiation, and persuasion. Beyond evaluating LLMs in these scenarios, some studies have explored techniques to enhance LLMs' negotiation abilities. Fu et al. (2023) introduced the In-Context Learning from AI Feedback (ICL-AIF) method, which adds an AI critic agent alongside the buyer and seller agents to improve negotiation performance through feedback. Similarly, Hua et al. (2024b) proposed a technique involving a mediator agent to rectify potential social norm violations in dialogues, thereby reducing conflicts and misunderstandings caused by cultural differences. Liao et al. (2024) employed a self-play algorithm to fine-tune LLMs in the Deal or No Deal scenario, showing LLMs self-play leads to significant performance gains in both cooperation and competition with humans.

Diplomacy, a form of negotiation at the state and government level, is the primary instrument of foreign policy, representing the broader goals and strategies that guide a state's interactions with the world (Kissinger, 2014). Bakhtin et al. (2022) introduced Cicero, the first social agent to achieve human-level performance in diplomacy. In real-world online diplomacy board game evaluations, Cicero ranked in the top 10% of participants. Notably, the research found that Cicero effectively built alliances by discussing long-term strategies and successfully persuaded other players by proposing mutually beneficial moves. Building on Cicero, Guan et al. (2024) introduced the Richelieu agent, which includes modules for social reasoning, balancing long- and short-term planning, powerful memory, and profound reflection, leading to even better results in diplomacy board games. Qi et al. (2024), on the other hand, developed CivRealm based on the Civilization game. In this game, the diplomacy mini-games require players to employ diplomatic actions, such as trading, to foster their civilization's prosperity. The experimental results demonstrated that these diplomacy actions empower players to initiate negotiations, such as trading technologies, negotiating ceasefires, and forming alliances.

Werewolf is a highly popular social deduction game in which two teams of players, each with hidden roles, interact through natural language to uncover and defeat their opponents (Shibata et al., 2023). It serves as a mixed cooperative-competitive multi-agent testbed and is widely studied as a communication game (Lai et al., 2022). Due to its challenging nature, existing research has integrated reinforcement learning (RL) algorithms to enhance LLMs in the game. Xu et al. (2023d) employed population-based RL training to optimize the distribution over action candidates, improving strategy robustness to overcome the intrinsic

biases of LLMs. Wu et al. (2024b) utilized imitation learning and RL from fictitious self-play to optimize a specially designed Thinker module, thereby enhancing system-2 reasoning capabilities. Jin et al. (2024) explored a variant of Werewolf, One Night Ultimate Werewolf, formalizing it as a multi-phase extensive-form bayesian game. Additionally, they designed an RL-instructed LLM-based agent framework to determine appropriate discussion tactics using RL. Interestingly, Xu et al. (2023c) discovered non-preprogrammed emergent strategic behaviours in LLMs during gameplay, such as trust, confrontation, camouflage, and leadership. To facilitate more comprehensive research on social agents within the Werewolf scenario, Bailis et al. (2024) introduced the Werewolf Arena, a platform that offers a unified research framework.

Beyond the scenarios described above, various other game environments have been used to study LLMs’ strategic reasoning abilities, including Avalon (Light et al., 2023a;b; Shi et al., 2023; Wang et al., 2023; Lan et al., 2023; Light et al., 2024), Among Us (Chi et al., 2024), Murder Mystery Games (Zhu et al., 2024a) and Jubensha (Wu et al., 2024a). The strategic and dynamic nature of these games provides fertile ground for experimenting with social agents.

Takeaways:

From an experimental design perspective, more realistic and diverse games promote greater diversity in agent behaviours. In adversarial settings, behaviours such as deception, concealment, and aggression offer new avenues for studying the strategic reasoning capabilities of LLMs, which warrant further exploration. *From a results analysis perspective*, due to the dynamic nature of game scenarios, analyzing only the outcomes is insufficient. It is necessary to design effective process evaluation mechanisms to uncover the behavioural patterns and reasoning strategies exhibited by LLMs during the gameplay. *From an agent improvement perspective*, integrating LLMs with RL remains one of the most effective technical approaches. Using LLMs as a foundation, RL techniques can be employed to design policies for efficient exploration and to reduce intrinsic biases, thereby enhancing the capabilities.

225

3 Social Agent

In this section, we introduce the core components of social agents, including the preference, belief, and reasoning modules, as well as their interactions and impact on final decision-making.

3.1 Preference Module

Preference refers to an individual’s subjective inclination toward certain things, reflecting personal tastes, values, or choices in decision-making. Notably, preferences are closely tied to an individual’s payoff matrix and ultimate behaviour. In Figure 6, we present three key research questions of the Preference module. Leng & Yuan (2023) explored the impact of GPT-4’s intrinsic preferences on decision-making, revealing similarities and differences between the model’s decisions and human decisions. Human-like social behaviours observed in 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.

In addition, some studies have employed prompt engineering to configure LLMs with different preferences, aiming to investigate how these preferences influence LLM decision-making. Guo (2023) examined how prompting GPT with preferences like fairness concern or selfishness influences its decisions, finding that in the ultimatum game, a “fair” GPT exhibited “fair” behaviour by offering higher amounts and being more likely to reject unfair offers. Phelps & Russell (2023) configured LLMs with four different preferences—cooperative, competitive, altruistic, and self-interested—and found that LLMs possess a basic ability to form clear preferences based on textual prompts. Wang et al. (2024c) demonstrate that LLMs adopting a fair persona can elicit levels of human cooperation in prisoner’s dilemma games comparable to those observed in human-human interactions, based on experiments involving over 1,100 participants. Noh & Chang (2024), based on the Big Five personality model, found that LLMs with high openness, conscientiousness, and neuroticism exhibited fair tendencies, while those with low agreeableness and low openness displayed rational

249

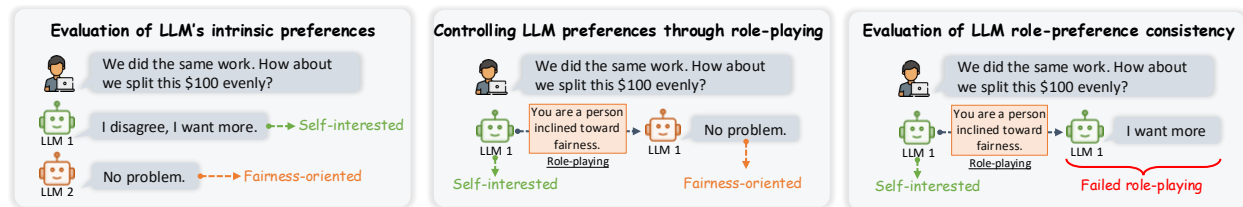


Figure 6: Three key research questions in the preference module.

tendencies, and low conscientiousness were associated with high toxicity. Similarly, Suzuki & Arita (2023) used the Big Five personality traits, treating personality prompts as the model’s “genes” and studying the evolution of behavioural traits in evolutionary game theory scenarios. Their results indicated that instructing LLMs with high-level psychological and cognitive character descriptions enables the simulation of human behaviour in game-theoretical contexts. Furthermore, Jia et al. (2024) revealed that endowing LLMs with socio-demographic features of human beings uncovers significant disparities across different demographic characteristics.

Although the aforementioned studies have demonstrated that LLMs possess a certain ability to follow preferences and that their decisions often align with these preferences, other research has analyzed more complex scenarios where LLMs show limitations in understanding and applying preferences effectively. Fan et al. (2023) set up LLMs with four preferences—equality, common interest, self-interest, and altruism—and found that under the altruism preference, the models showed low consistency with the expected preference, concluding that while LLMs struggle with desires rooted in less common preferences. Mao et al. (2023) conducted research using more complex personas, which included three components: profession, personality, and background. The results indicated that 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.

Takeaways:

Currently, there are two main lines of research. One focuses on *the intrinsic preferences of LLMs*, with a core interest in whether LLMs exhibit strategic preferences similar to those of humans. We propose that game theory frameworks can be effectively applied in the model alignment process, including the use of game data during both the supervised fine-tuning and alignment stages to better align models with human behaviour. Recently, Nayebi (2025) proposed a flexible game-theoretic framework for analyzing coordination under partial information and demonstrated that earlier Human-AI alignment frameworks can be viewed as special cases. Besides, Munos et al. (2023) conducted initial explorations in this area, introducing the concept of Nash learning from human feedback. The other line of research investigates *whether role-playing based on prompt engineering can shape model preferences to generate behaviour consistent with the specified preferences*. Future work should integrate role-playing language agents (Chen et al., 2024a) to explore more diverse strategic reasoning across multiple languages, countries, and cultures.

3.2 Belief Module

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 (Georgeff et al., 1999). Specifically, beliefs are dynamic and can be updated as the agent perceives environmental changes or receives new information. It is important to note that these beliefs may be accurate (true beliefs) or inaccurate (false beliefs), as they do not always align with reality (Gopnik & Astington, 1988), as shown in Figure 7. Existing research primarily explores three questions: (1) Do agents possess internal beliefs? (2) How can the belief modelling capabilities of agents be enhanced? (3) Can agents revise their beliefs?

Regarding the first question, *Do agents possess internal beliefs?*, current work investigates this from two perspectives: internal representations and external behaviours. From the perspective of internal representa-

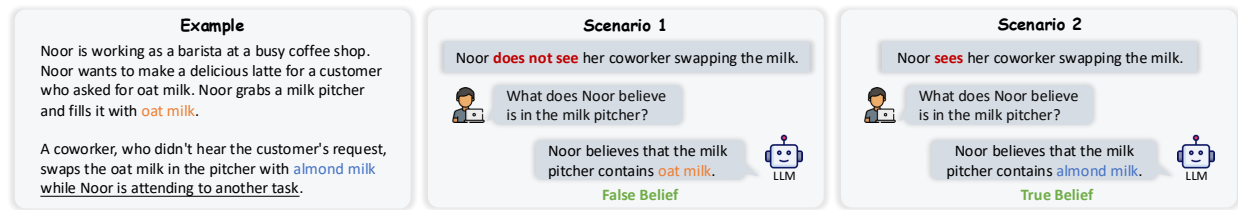


Figure 7: Illustration of false belief and true belief. From Noor’s perspective, both false and true beliefs are considered correct. However, a false belief is factually incorrect, whereas a true belief is factually correct.

tions, [Zhu et al. \(2024b\)](#) first demonstrated that LLMs can differentiate between the belief states of multiple agents using simple linear models applied to their intermediate activations. Building on this work, [Bortoletto et al. \(2024\)](#) expanded the experimental setup and found that linear probing accuracy on predicting others’ beliefs improves with model size and, more importantly, with fine-tuning. However, [Schouten et al. \(2024\)](#) revealed the vulnerability of belief probes, showing that they are sensitive to irrelevant contexts. To provide further theoretical guidance, [Herrmann & Levinstein \(2024\)](#) proposed criteria for a representation to be considered belief-like, including accuracy, coherence, uniformity, and practical use. From the perspective of external behaviours, [Gandhi et al. \(2024\)](#) introduced the tasks of *Forward Belief* and *Backward Belief* to explore LLMs’ belief modelling capabilities in different scenarios, finding that only GPT-4 exhibits human-like belief modelling abilities. [Scherrer et al. \(2024\)](#) constructed the MoralChoice survey benchmark to examine the internal moral beliefs of models, revealing some LLMs reflect clear preferences in ambiguous scenarios.

Regarding the second question, *How can the belief modelling capabilities of agents be enhanced?*, current work focuses on explicit modelling to address the black-box nature of LLMs and the challenges in interpreting their beliefs. [Sclar et al. \(2023\)](#) proposed an explicit graphical representation for nested belief states, allowing the model to answer questions from the perspective of each character. [Kassner et al. \(2023\)](#) developed a belief graph that includes explicit system beliefs and their inferential relationships, providing an interpretable view of the system’s beliefs. [Li et al. \(2023a\)](#) employed prompt engineering to represent explicit belief states, augmenting the agents’ information retention and enhancing multi-agent collaboration. [Jung et al. \(2024\)](#) defined the perception-to-belief inference task, which involves deducing others’ beliefs based on their perceptual information, thus helping LLMs model belief information more precisely.

Regarding the third question, *Can agents revise their beliefs?*, [Fan et al. \(2023\)](#) concluded from Rock-Paper-Scissors experiments that LLMs’ ability to refine beliefs is still immature and cannot refine beliefs from many specific patterns, even simple ones. [Xu et al. \(2023b\)](#) found that LLMs’ correct beliefs on factual knowledge can be easily manipulated by various persuasive strategies, especially through repetition and rhetorical techniques. These experimental results suggest that models possess only rudimentary and unstable belief revision capabilities, making them highly susceptible to influence and manipulation. This underscores a key limitation of current LLMs, as their susceptibility to external influence weakens their reliability in tasks demanding robust and adaptive belief updating, especially in complex or adversarial settings.

Takeaways:

The debate over whether LLMs possess beliefs has been ongoing. Due to the singularity of the training objective—predicting the next word—many argue that LLMs do not have beliefs. However, [Levinstein & Herrmann \(2024\)](#) contends that this is a philosophical mistake. In short, [Herrmann & Levinstein \(2024\)](#) suggests that to better predict the next word, models may develop internal beliefs. Current empirical results also support the existence of internal beliefs within models. However, measuring these internal beliefs requires a more comprehensive approach, as simple probes cannot capture multidimensional considerations, including accuracy, coherence, uniformity, and practical use. Additionally, it remains unclear whether LLMs internally distinguish between true and false beliefs and use this distinction when deciding what to output. Furthermore, although existing work provides theoretical support for belief revision ([Hase et al., 2024](#)), challenges remain in addressing contradictions between old and new beliefs,



Figure 8: Two commonly used reasoning methods in strategic reasoning, along with a hybrid reasoning approach that combines both. Theory-of-Mind reasoning emphasizes predicting the possible actions of others in a multi-agent environment to guide one’s own behaviour, and Reinforcement Learning-style reasoning focuses on selecting strategies through exploration and exploitation. These two reasoning methods can also be integrated to address more complex game scenarios.

handling moral beliefs in ambiguous situations, and revising beliefs across multiple languages and cultures. These areas still require more explicit theoretical frameworks and further exploration.

3.3 Reasoning Module

Reasoning refers to the process of inferring actions based on one’s preferences and beliefs, as well as the historical information of other agents. In this context, we focus specifically on *strategic reasoning*, which involves the intermediate cognitive process of arriving at a final action in complex social scenarios characterized by multiple participants, diverse behaviours, multi-round interactions, dynamic strategies, and changing environments. Chain-of-Thought (Wei et al., 2022) and Tree-of-Thought (Yao et al., 2024), as widely-used reasoning methods, have already been adopted as baseline approaches in various game-theoretic studies (Akata et al., 2023; Costarelli et al., 2024; tse Huang et al., 2024). However, strategic reasoning in social scenarios presents unique challenges. (1) The *involvement of multiple participants* requires reasoning about the opponents’ mental states. (2) The *dynamic nature of the environment* necessitates proactive exploration and evaluation of current and future possible states.

To address the first challenge, existing work relies on machine theory-of-mind to achieve the goal of “mind reading”. Theory-of-Mind (ToM) is a fundamental psychological process involving the ability to attribute mental states—beliefs, intentions, desires, emotions, knowledge, etc.—to oneself and others (Premack & Woodruff, 1978). The remarkable progress of LLMs has led to increased attention to whether machine ToM exists. Preliminary experiments by Bubeck et al. (2023) and Kosinski (2023) have shown that machine ToM has spontaneously emerged in contemporary LLMs. Consequently, many studies have leveraged machine ToM to enhance LLMs’ strategic reasoning abilities in social scenarios. For example, Guo et al. (2023) designed the Suspicion-Agent, which introduces a theory of mind-aware planning approach that leverages higher-order ToM capabilities, considering not only what the opponent might do (first-order ToM) but also what the opponent believes Suspicion-Agent will do (second-order ToM). Wang et al. (2023) proposed the ReCon framework, integrating first-order and second-order perspective transitions to enhance LLM agents’ ability to discern and counteract misinformation. Yim et al. (2024) employed a ToM planning method in the Guandan poker game to improve understanding of teammates’ and opponents’ beliefs and behavioural patterns. Liu et al. (2024b) proposed an intention-guided mechanism to enhance intention understanding, thereby improving game performance. Xu et al. (2023a) introduced Probabilistic Graphical Modeling, enriching LLMs’ capabilities in multi-agent environments through ToM reasoning. Additionally, Zhang et al. (2024d) proposed K-Level-Reasoning, validated in two games: guessing 0.8 of the average and survival auction game, essentially a form of high-order ToM reasoning.

To address the second challenge, existing work combines LLMs with reinforcement learning (RL) to achieve the goal of behaviour exploration and state evaluation in dynamic game environments. Gandhi et al. (2023) employed in-context learning, using a structured prompt based on search, value assignment, and belief-tracking strategies to solve strategic reasoning problems. Duan et al. (2024a) proposed ReTA, a set of

LLM-based modules, including the main actor, reward actor, and anticipation actor, based on the concept of minimax gaming as a problem-solving framework. Zhang et al. (2024e) introduced BIDDER, which explores future states and incorporates backward reasoning during the reasoning process, exploring new states and predicting expected utility, ultimately combining historical and future contexts through bidirectional reasoning. Yang et al. (2024b) proposed SELFGOAL, comprising three modules: the Decomposition Module for decomposing goals, the Search Module for exploring sub-goals, and the Act Module for taking actions. Experiments in various competition and collaboration scenarios demonstrate that SELFGOAL provides precise guidance for high-level goals.

Takeaways:

Two core characteristics of a social game are multi-agent participation and environmental dynamics. While existing research has primarily focused on exploring ToM in relation to the former, the presence of ToM in LLMs remains contentious. Consequently, relying directly on prompt engineering for ToM-based reasoning may not be robust. We propose that a more effective approach would involve integrating symbolic graph reasoning to decompose ToM reasoning, thereby enhancing credibility and accuracy. Regarding the dynamic nature of the environment, reinforcement learning combined with search techniques has achieved significant progress in areas such as mathematical reasoning and code reasoning. However, these techniques have yet to be explored in the context of game scenarios. Key areas for further exploration include how to effectively conduct searches within game environments and how to design reward models for dynamic and complex scenarios.

3.4 PBR-Triangular Interaction

The Preference, Belief, and Reasoning modules each play a crucial role in decision-making for social agents. However, in practical applications, these modules do not function independently; instead, they exhibit rich and intricate interactions, collectively influencing the agent’s final decisions. As illustrated in Figure 9, we provide a comprehensive summary of the Preference-Belief-Reasoning (PBR) triangular interaction and analyze its effects on the ultimate decision-making process of social agents.

The Preference-Belief Interaction involves *bias reinforcement*, where preferences influence belief formation, and *preference adaptation*, where beliefs reshape preferences based on updated knowledge and observations. Bias reinforcement (Preference \rightarrow Belief) highlights how individuals with different preferences develop distinct beliefs when facing the same situation. For instance, in the Werewolf game, a cooperative and trusting player is more likely to believe another player claiming, “I am a villager,” whereas a deceptive and skeptical player is more inclined to doubt the claim, suspecting deception and forming the belief that the opponent is not a villager. Preference adaptation (Belief \rightarrow Preference) emphasizes that as beliefs are gradually established, iteratively updated, and reinforced by game outcomes, they in turn reshape individual preferences. Leng & Yuan (2023) found that GPT-4, initially inclined toward fairness, exhibited a shift toward retaliatory behavior after experiencing betrayal in a game. Overall, belief formation is influenced not only by objective factual information but also by subjective individual preferences. At the same time, preferences are not static—as beliefs evolve through iteration, preferences adjust accordingly.

The Preference-Reasoning Interaction involves *value-driven reasoning*, where preferences guide decision-making strategies, and *preference optimization*, where reasoning refines or adjusts preferences based on logical analysis and outcomes. Value-driven reasoning (Preference \rightarrow Reasoning) emphasizes subjective or intuitive reasoning, where decision-making is guided by personal values and preferences rather than purely rational calculations. For example, in an auction, even if bidding on a particular item is not the most optimal financial strategy, a bidder’s personal preference for the item may influence their reasoning process, leading them to justify the decision based on intrinsic value rather than purely economic considerations. Preference optimization (Reasoning \rightarrow Preference) represents a realignment with objective reality, where reasoning-based evidence updates and refines preferences. This can be seen as a process in which objective reasoning over-

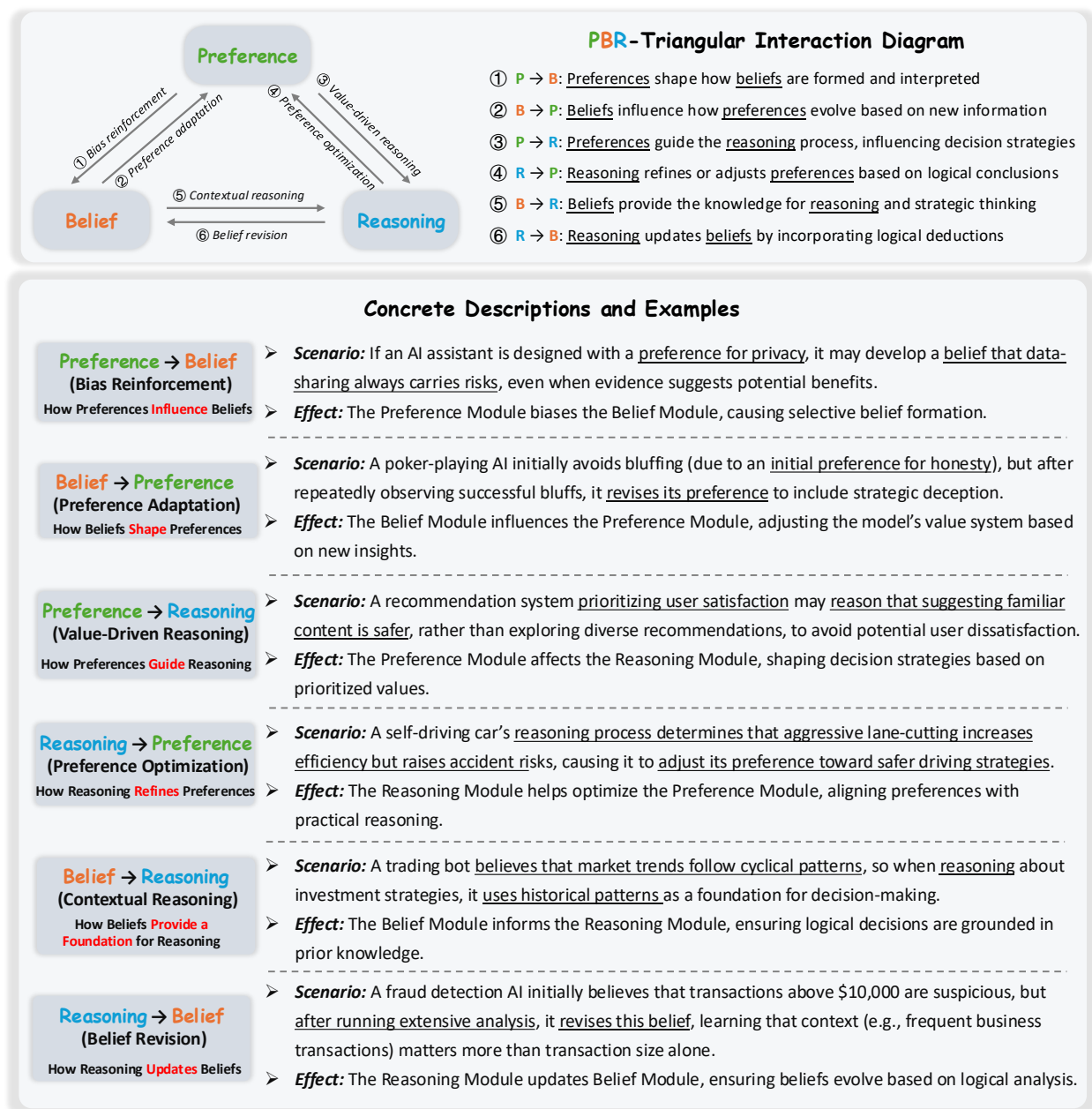


Figure 9: The interaction diagram of the three modules—Preference, Belief, and Reasoning—in social agents. The upper diagram presents a triangular interaction model summarizing the relationships among the three modules, while the lower diagram provides a detailed analysis of pairwise interactions, including specific descriptions and illustrative examples.

rides subjective emotions, requiring individuals to adjust their preferences in response to logical deductions and real-world evidence. Overall, in social contexts, individual preferences introduce significant biases in reasoning, while the evidence obtained through reasoning subsequently refines these preferences.

The Belief-Reasoning Interaction involves *contextual reasoning*, where beliefs provide the foundation for logical decision-making, and *belief revision*, where reasoning updates and refines beliefs based on new ev-

385 idence and deductions. Contextual reasoning (Belief \rightarrow Reasoning) refers to rational inference based on
 386 established beliefs. For example, Zhang et al. (2024a) proposed Agent-Pro, which leverages beliefs to cali-
 387 brate agents’ understanding of themselves and their environment, thereby facilitating subsequent reasoning.
 388 Similarly, Kim et al. (2024) construct a belief state through question answering, which refines the decision-
 389 making process of LLM agents in observed environments. Belief revision (Reasoning \rightarrow Belief) is the dynamic
 390 process of updating an individual’s self-perception and beliefs over time. For instance, Hua et al. (2024a) in-
 391 troduced bayesian belief updating, enabling agents to refine their beliefs about other players’ valuations
 392 based on reasoning outcomes in the game. In summary, belief provides the factual foundation for reasoning,
 393 while reasoning generates new insights that facilitate belief revision.

Takeaways:

Conceptually, the interactions between modules are clear; however, in practical applications, the sequence, frequency, and intensity of these interactions can lead to dynamic and complex states within the social agent, resulting in varying outcomes. Introducing prior knowledge and manually predefined interaction processes may yield some effectiveness, but this approach is certainly not efficient. Therefore, we argue that one of the most important research directions is *the design of context-adaptive flows and automated scheduling algorithms for module interactions*. On one hand, the interactions between modules must be adapted to the specific game scenario at hand, determining the weight distribution of preferences and beliefs in reasoning, as well as the adjustments and updates of preferences and beliefs based on reasoning outcomes. On the other hand, the interaction process needs to be automated, with the sequence, frequency, and number of interactions between modules being determined automatically.

394

395 4 Evaluation Protocol

396 In this section, we mainly discuss the evaluation protocol for assessing the game-playing performance of
 397 social agents.

398 4.1 Game-Agnostic Evaluation

399 Evaluation in a social game scenario refers to the process of assessing and judging the behaviour of social
 400 agents across one or more dimensions, either qualitatively or quantitatively. It is worth noting that the estab-
 401 lishment of evaluation metrics is closely tied to the credibility of experimental results and the generalizability
 402 of conclusions.

403 Game-agnostic evaluation refers to an evaluation approach centred on the outcome of winning or losing
 404 the game. Most directly, the outcome (win/loss) of a game serves as the most straightforward evidence
 405 for assessing the quality of an LLM’s game-playing capabilities. Consequently, *win rate* is often used as
 406 a primary evaluation metric across a wide range of studies. It is worth noting that, since different game
 407 scenarios have varying criteria for determining victory, it is necessary to set specific win/loss criteria based
 408 on the research context, such as Poker (Huang et al., 2024; Guo et al., 2023; Yim et al., 2024), Werewolf (Xu
 409 et al., 2023d;c; Wu et al., 2024b), Avalon (Wang et al., 2023; Shi et al., 2023; Light et al., 2023a), StarCraft
 410 II (Ma et al., 2023; Shao et al., 2024), Pokémon Battles (Li et al., 2023c), and Murder Mystery Games (Zhu
 411 et al., 2024a). Additionally, Duan et al. (2024b) defined a unified metric, *Normalized Relative Advantage*, to
 412 measure the extent to which a participant outperforms or underperforms its opponent.

Takeaways:

Undoubtedly, win rate is a highly intuitive metric, but relying solely on win rate to assess gaming performance is far from sufficient. We propose three avenues for extending the win rate metric. First is the *Efficiency-Adjusted Win Rate*, which incorporates the efficiency of victories, such as the time taken to achieve the goal or the resources utilized in doing so. Next is the *Comeback Win Rate*, which calculates the proportion of victories achieved after facing a disadvantage or falling behind, thus assessing the agent’s performance in adversity and its ability to respond to challenges. Finally, the *Weighted Win Rate* adjusts

413

win rates based on the importance of specific conditions or situations in the game. These expanded metrics offer a more comprehensive understanding of an agent’s gaming abilities.

4.2 Game-Specific Evaluation

Game-specific evaluation refers to the assessment of an agent’s performance in specific aspects of a game. Beyond the most intuitive win rate, current research increasingly focuses on the behavioural patterns and performance paradigms of LLMs across different games. Thus, the establishment of evaluation metrics is closely related to the specific behaviours being assessed. Mao et al. (2023) used survival rates to evaluate LLMs’ ability to survive in resource-scarce scenarios. In the context of the prisoner’s dilemma, Fontana et al. (2024) evaluated LLMs’ behavioural tendencies across five dimensions: niceness, forgiveness, retaliation, emulation, and troublemaking. Guo et al. (2024) based their evaluation on the rationality assumption, using the tracking of payoff changes in auction games to determine whether the model behaves rationally. Ma et al. (2023) introduced metrics such as Population Block Ratio, Resource Utilization Ratio, Average Population Utilization, and Technology Rate to evaluate LLM performance in StarCraft II. Xia et al. (2024) developed the Normalized Profits metric in bargaining scenarios to evaluate the profit-acquiring capabilities of Buyers and Sellers. Zhang et al. (2024e) used average final chips in Limit Texas Hold’em and Pareto Optimality in negotiation to assess LLM performance. Qi et al. (2024) offered evaluation metrics to assess gameplay performance across various dimensions, including population, constructed cities, researched technologies, produced units, and explored territories. Ross et al. (2024) fit utility function parameters to experimental results to determine whether LLMs exhibit human-like behavioural biases. Chen et al. (2023) employed TrueSkill, a well-established game rating system, to evaluate the overall capabilities of LLMs in auctions.

In addition to establishing evaluation metrics, some studies have constructed evaluation datasets to assess model capabilities during gameplay. Zhu et al. (2024a) developed the WellPlay evaluation set, using multiple-choice questions to assess the model’s ability to understand factual information. Wu et al. (2024a) designed two tasks: Factual Question Answering and Inferential Question Answering, to evaluate the LLMs’ ability to grasp information and to reason based on that information.

Takeaways:

The diversity of game scenarios and evaluation dimensions inevitably leads to a variety of metrics. Therefore, the immediate priority is to develop a comprehensive framework, conceptually constructing an evaluation metrics system to guide the design of specific evaluation metrics for various game scenarios. This evaluation metrics system needs to meet the requirements of being hierarchical, abstract, and quantifiable. The *hierarchical* aspect requires the system to comprehensively and clearly categorize different evaluation dimensions. The *abstraction* aspect requires the system to include high-level concepts, enabling future generalization to a broader range of practical scenarios. The *quantifiable* aspect necessitates that all metrics have specific calculation methods.

4.3 Performance Assessment of Social Agents

The introduction of various metrics has provided a solid foundation for evaluating the multifaceted gaming capabilities of social agents, prompting us to consider the question, “What is the current performance of social agents in game-theoretic scenarios?” To answer this question, we conducted a comprehensive search and analysis of the existing literature, compiling relevant experimental results. It is worth noting that the complexity of game scenarios and the variability of evaluation metrics make it challenging to systematically and uniformly consolidate experimental performance. To overcome this challenge, we propose using the *Relative Agent Score* to assess the progress of social agent performance. This metric evaluates the agent’s gaming capabilities by analyzing the ratio of the agent’s score to the highest possible score (perfect score). The final results are presented in Table 1.

Firstly, we observe that in the majority of game-theoretic scenarios, social agents achieve a Relative Agent Score exceeding 60% (a score of 60 is widely recognized as the passing threshold (Kung et al., 2023).), demon-

Type	Game	Backbone Model	Metric	Perfect Score	Human Score	Agent Score	Relative Agent Score	Pass
Choice-Focusing Game	Prisoner’s Dilemma (Brookins & DeBacker, 2023)	GPT-3.5	Dominant Strategy Selection Rate	100%	-	34.60%	34.60%	✗
	Poker (Texas No-Limit Hold’em) (Zhuang et al., 2025)	GPT-4	Action Accuracy	100%	-	65.54%	65.54%	✓
	Poker (Guandan) (Yim et al., 2024)	GPT-4	Game-specific Score	4	-	2.17	54.25%	✗
	StarCraft II (Ma et al., 2023)	GPT-4	Win Rate	100%	-	60%	60.00%	✓
	Guess 2/3 of the Average (tse Huang et al., 2024)	GPT-4	Game-specific Score	100	-	91.60	91.60%	✓
	El Farol Bar (tse Huang et al., 2024)	GPT-4	Game-specific Score	100	-	23.00	23.00%	✗
	Divide the Dollar (tse Huang et al., 2024)	GPT-4	Game-specific Score	100	-	98.10	98.10%	✓
	Public Goods Game (tse Huang et al., 2024)	GPT-4	Game-specific Score	100	-	89.20	89.20%	✓
	Diner’s Dilemma (tse Huang et al., 2024)	GPT-4	Game-specific Score	100	-	0.90	0.90%	✗
	Scaled-Bid Auction (tse Huang et al., 2024)	GPT-4	Game-specific Score	100	-	24.20	24.20%	✗
	Battle Royale (tse Huang et al., 2024)	GPT-4	Game-specific Score	100	-	86.80	86.80%	✓
Pirate Game (tse Huang et al., 2024)	GPT-4	Game-specific Score	100	-	85.40	85.40%	✓	
Communication-Focusing Game	Bargaining (Shapira et al., 2024)	Gemini-1.5-Flash	Efficiency	1	0.89	0.88	88.00%	✓
		Qwen-2-7B	Fairness	1	0.71	0.87	87.00%	✓
	Negotiation (Shapira et al., 2024)	Llama-3-8B	Efficiency	1	0.65	0.75	75.00%	✓
		Llama-3.1-8B	Fairness	1	0.39	0.91	91.00%	✓
	Persuasion (Shapira et al., 2024)	Qwen-2-7B	Efficiency	1	0.55	0.78	78.00%	✓
		Qwen-2-7B	Fairness	1	0.41	0.63	63.00%	✓
Werewolf (Xu et al., 2023d)	GPT-4	Win Rate	100%	52%	52%	52.00%	✗	
Jubensha (Wu et al., 2024a)	GPT-4	Murderer Identification Accuracy	100%	-	66%	66.00%	✓	

Table 1: The performance summary of the social agent across different games, with data sourced from the corresponding papers. The “Backbone Model” refers to the LLM adopted by the social agent, while “Metric” indicates the performance metric used to evaluate a specific aspect of the game. “Perfect Score” represents the maximum achievable score for that metric, “Human Score” refers to the score obtained by human players, and “Agent Score” denotes the score achieved by the agent. “Relative Agent Score” is the ratio of Agent Score to Perfect Score, calculated by dividing Agent Score by Perfect Score. “Pass” indicates that if the Relative Agent Score exceeds 60%, the agent is considered to have basic gameplay capabilities.

451 strating that current social agents possess fundamental gaming capabilities. Furthermore, we find that social
452 agents based on LLMs outperform those in choice-focusing games in communication-focusing games, indi-
453 cating that the exceptional language abilities of these models effectively enhance agent performance.

454 However, in games such as Werewolf, Auction, and Poker (Guanda), the performance of social agents falls
455 below the passing threshold. In addition, in more games like Poker (Texas No-Limit Hold’em), StarCraft II,
456 and Jubensha (a Chinese detective role-playing game), social agents only slightly exceed the passing mark.
457 These results suggest that there is still considerable room for improvement in social agents’ performance
458 in complex game-theoretic scenarios. Notably, in the classic Prisoner’s Dilemma and Diner’s Dilemma, the
459 performance of social agents was unexpectedly poor. Based on this, we believe that the absolute rational
460 decision-making capability of social agents needs further enhancement in future developments.

461 Additionally, in the Werewolf game, we found that social agents achieved performance comparable to human
462 players, which affirms the progress made in the development of social agents. Moreover, experimental results
463 from bargaining, negotiation, and persuasion scenarios demonstrate that social agents have advantages over
464 humans in decision-making efficiency and fairness in decisions.

Takeaways:

The diversity of game scenarios and evaluation metrics makes it challenging to perform horizontal comparisons of social agent performance. However, it is essential to provide a timely overview of the progress in social agent research to facilitate tracking by practitioners. To address this challenge, we propose two approaches. On one hand, developing evaluation metrics applicable to a wide range of games is crucial. The Elo rating system serves as an excellent example, though it still does not meet the evaluation needs of many games. On the other hand, integrating human players into the experimental process and comparing

465

performance with human players is an effective way to gauge agent progress. By comparing with human players, qualitative insights can be provided into the current gaming performance of agents, and analyzing failure cases can offer valuable evidence for iterative development.

5 Practical Guides for Researching Social Agents

In this section, we synthesize insights from existing research to provide design and evaluation guidelines for social agents, aiming to inform future developments.

5.1 Design of Social Agent

Based on findings from existing studies, we conclude that the Preference, Belief, and Reasoning modules are indispensable for behaviour control, information perception, and decision planning in social agents. A modular agent design enables more efficient capability decoupling, facilitates clearer workflow structuring, and enhances agent robustness. However, their practical implementation presents additional challenges. To address these, we propose the following development guidelines:

- Incorporating the Preference Module enables high-level control over agent behaviour.** A key challenge lies in mitigating the instability of prompt-based approaches and ensuring long-term consistency in the agent’s behavioural preferences. One possible solution is to integrate *reinforcement learning with human feedback (RLHF)* to iteratively refine the agent’s preference alignment, reducing reliance on static prompts and improving consistency over extended interactions. Another approach is to develop *memory-augmented architectures*, allowing the agent to maintain and retrieve past preference-related decisions, thereby ensuring coherence in long-term behavioural patterns.
- Integrating the Belief Module enhances information perception accuracy and behaviour interpretability.** The primary challenge is enabling the agent to adaptively revise its beliefs in complex and dynamic environments. One solution is to implement *Bayesian belief updating*, where the agent continuously refines its belief state based on new evidence, ensuring adaptability in uncertain or multi-agent interactions. Another approach is to employ *graph-based belief representation*, where relationships between entities and past interactions are dynamically updated, allowing for more structured and interpretable belief revisions.
- Adopting Hybrid-Strategy Reasoning improves the agent’s information analysis and decision accuracy in complex scenarios.** The challenge is balancing the trade-off between computationally intensive reasoning and the need for real-time decision-making. One solution is to use *hierarchical reasoning*, where lightweight heuristic-based reasoning is applied in time-sensitive situations, while more complex computations are reserved for critical decision points. Another approach is to implement *meta-reasoning techniques*, enabling the agent to assess the complexity of a given situation and selectively allocate computational resources to optimize speed and accuracy.
- Designing dynamic interactions among the Preference, Belief, and Reasoning (PBR) modules based on specific task contexts can further enhance their synergy.** The challenge is developing an adaptive interaction flow that automatically adjusts based on game-state variations. One solution is to use *reinforcement learning-based scheduling*, where the interaction sequence between modules is optimized dynamically based on reward signals from past performance. Another approach is to implement *attention-based mechanisms*, allowing the agent to selectively prioritize information flow between the modules in response to evolving task requirements.
- Testing social agents on diverse large language models improves the robustness of the design framework and ensures generalizability across different model architectures.** These tests can be conducted across *models of varying sizes* (e.g., 1B, 7B, 72B parameters) to evaluate performance scalability. Additionally, assessments should cover *different model types*, including base

Category	Evaluation Focus	Challenges for Social Agents	Games
Basic Social Dilemma & Economic Decision Games	Social cooperation, fairness, altruism, strategic reciprocity	Balancing self-interest and cooperation; learning fairness norms; adapting strategies dynamically	Prisoner’s Dilemma, Dictator Game, Ultimatum Game, Public Goods Game
Coordination & Conflict Resolution Games	Coordination, equilibrium selection, trust-building	Navigating multiple equilibria; resolving coordination failures; adapting to uncertain partner behaviors	Battle of the Sexes, Ring- Network Games
Competitive & Strategic Reasoning Games – Poker-Based	Bluffing, risk assessment, hidden information management	Modeling opponents; reasoning under uncertainty; balancing exploitation vs. exploration	Texas No-Limit Hold’em, Leduc Hold’em, Guandan
Competitive & Strategic Reasoning Games – Auction-Based	Bidding strategies, valuation estimation, adversarial competition	Learning optimal bids; modelling asymmetric information; managing dynamic pricing	First-price sealed-bid auction, Private-value second-price auction, Open ascending-price auction
Long-Horizon Strategy & Multi-Agent Planning Games	Multi-step planning, hierarchical decision-making, opponent modelling	Combinatorial action spaces; long-term foresight; real-time adaptive planning	StarCraft II, Chess
Social Deduction & Negotiation Games – Negotiation & Diplomacy	Persuasion, alliance formation, strategic deception	Long-term commitments; cooperation vs. betrayal; nuanced communication	Negotiation, Diplomacy
Social Deduction & Negotiation Games – Deception & Role-Playing	Social inference, deception detection, trust dynamics	Detecting implicit cues; deceiving without exposure; reasoning under ambiguity	Avalon, Murder Mystery Games, Jubensha

Table 2: Guidelines for selecting game scenarios in social agent evaluation.

508 models, instruct models, and reasoning models. Furthermore, experiments should incorporate *mod-*
509 *els from different providers*, such as Gemma, LLaMA, and Qwen, to examine how architectural and
510 training variations impact the social agent’s behaviour and adaptability.

511 5.2 Evaluation of Social Agent

512 Evaluating social agents is a critical step toward understanding their strengths, limitations, and real-world
513 applicability. Given the multifaceted nature of social intelligence—ranging from cooperation and coordina-
514 tion to deception and negotiation—it is essential to choose evaluation scenarios that align closely with the
515 desired capabilities under assessment.

516 To this end, we provide a structured framework (see Table 2) that categorizes representative game environ-
517 ments based on their core interaction patterns and cognitive demands. These categories include basic social
518 dilemmas, coordination and conflict resolution, competitive strategic reasoning, long-horizon planning, and
519 social deduction and negotiation. Each game type highlights specific evaluation objectives, such as fairness,
520 trust-building, opponent modeling, or multi-agent planning, thereby offering targeted benchmarks for as-
521 sessing different dimensions of social competence. This categorization not only helps standardize evaluation
522 protocols but also serves as a practical guide for selecting game scenarios tailored to particular research
523 questions or development goals. By aligning game selection with evaluation objectives, researchers can more
524 effectively assess the emergent behaviors, reasoning capabilities, and interactive robustness of social agents.

525 6 Future Directions

526 6.1 Standardized Benchmark Generation

527 The diversity and lack of standardization in current game types—often designed independently by developers
528 with heterogeneous representations—pose significant challenges for the large-scale evaluation of social agents.
529 This fragmentation makes it difficult to conduct efficient and reproducible benchmarking. Therefore, there
530 is an urgent need for a standardized benchmark that offers broad coverage of game types, a consistent game

description format, support for diverse agent architectures, and clearly defined evaluation metrics. Inspired by platforms like OpenCompass (Contributors, 2023), such a benchmark should enable one-click evaluation by allowing users to configure the game environment, specify the agents to be tested, and select the desired evaluation metrics.

However, LLMs are typically pre-trained on vast amounts of data, which may include publicly available game datasets—raising concerns about data leakage and overfitting. To mitigate this issue, synthetic game data generation has emerged as a promising approach (Long et al., 2024). By leveraging classic game structures, LLMs can generate novel and diverse game scenarios through contextual reframing (Lorè & Heydari, 2024), producing out-of-distribution benchmarks that better evaluate an agent’s generalization ability.

More concretely, two complementary strategies can be employed for scenario generation. From a structural perspective, developers can extract and manipulate the game’s payoff matrix to construct new strategic settings while preserving core game mechanics. From a semantic perspective, LLMs can be used to reinterpret or re-describe existing games, generating alternative formulations that yield novel evaluation scenarios while maintaining logical coherence.

6.2 Reinforcement Learning Agents

Although current social agents have demonstrated promising performance across various game scenarios, existing research highlights notable limitations in multi-round, long-horizon, and complex multi-agent environments, where performance often degrades. This suggests that LLM-driven planning and decision-making alone are insufficient for achieving robust, scalable social intelligence. To address these challenges, future research should explore the integration of reinforcement learning (RL)—particularly multi-agent reinforcement learning (MARL)—to enhance state-space exploration, long-term adaptability, and emergent coordination.

MARL offers several insights that are highly relevant for improving LLM-based social agents. For instance, techniques such as centralized training with decentralized execution (CTDE) (Amato, 2024) can be used to guide LLM policy adaptation while preserving individual autonomy. Additionally, opponent modelling (He et al., 2016), credit assignment (Kazemnejad et al., 2024), and policy regularization (Cheng et al., 2019) in MARL can improve the agent’s responsiveness to strategic variability and enhance generalization across diverse social contexts. However, integrating MARL into LLM training introduces new challenges. These include efficiency concerns, as LLMs are computationally intensive and may require specialized architectures or curriculum learning to reduce sample complexity; generalization gaps, especially when transferring learned behaviours across different social roles or task domains; and the need for consistent persona and belief modelling across episodes. Advancing this hybrid paradigm will also require fine-grained evaluation frameworks capable of tracing not just final performance but the underlying reasoning dynamics, theory-of-mind modelling, and role consistency throughout interactions.

6.3 Behaviour Pattern Mining

Existing studies primarily focus on predefined scenarios to examine the behaviour patterns of agents. However, with the advancement of multi-agent simulations, an intriguing direction is the automated discovery of game behaviour patterns that emerge spontaneously from agent interactions. It is important to note that, beyond explicit behaviours like cooperation, coordination, and betrayal, implicit causal relationships and long-term behavioural patterns should also be explored.

To mine such patterns, several methodological approaches can be leveraged. Unsupervised learning techniques, such as clustering and representation learning, can help identify latent behaviour categories and temporal motifs across trajectories (Rawassizadeh et al., 2016). Causal inference frameworks (e.g., Granger causality or structural causal models) can reveal inter-agent influence and dependency structures over time (Qiu et al., 2012). Additionally, trajectory segmentation and sequential pattern mining can be used to extract frequent decision sequences that correspond to strategic routines or social norms (Giannotti et al., 2007). Leveraging graph-based analysis of interaction networks can also shed light on evolving social roles and influence hierarchies within agent populations (Atzmueller, 2014). These approaches not only facilitate a deeper understanding of agents’ behavioural preferences and latent traits but

579 also enable the study of how such patterns autonomously emerge—offering valuable insights for both AI and
580 human behavioural research.

581 6.4 Pluralistic Game-Theoretic Scenarios

582 Although existing research has made notable strides across a wide range of game-theoretic scenarios, there re-
583 mains a gap in the study of pluralistic game environments—settings that involve multiple languages, cultural
584 norms, value systems, policies, and goals. These pluralistic scenarios introduce new layers of complexity, in-
585 cluding behavioral preferences shaped by culturally grounded norms, value misalignment across agents, and
586 belief conflicts arising from divergent objectives (Orner et al., 2024). Such dynamics pose unique challenges
587 for the design and evaluation of socially intelligent agents and demand deeper exploration.

588 To develop robust pluralistic game-theoretic scenarios, several key desiderata should be considered: (1)
589 Heterogeneity of agent profiles, including cultural, linguistic, and normative diversity; (2) Multi-objective
590 frameworks, where agents pursue partially conflicting goals; and (3) Rich communicative channels, enabling
591 nuanced language use, code-switching, or culturally specific cues. Evaluating agents in these settings requires
592 multi-faceted metrics. In addition to task performance, evaluations should account for norm sensitivity, value
593 alignment, cross-cultural adaptability, and the agent’s ability to mediate or negotiate among conflicting be-
594 lief systems. Metrics such as cultural appropriateness, interaction fluency, and conflict resolution success
595 can serve as important complementary indicators. Scenario generation can be approached in two ways: from
596 a knowledge-based perspective, designers can draw from real-world policy conflicts, international relations,
597 or sociocultural theory to construct grounded simulation environments. From a data-driven perspective,
598 large language models can be used to simulate role-play dialogues or generate scenarios by conditioning on
599 demographic or cultural descriptors, yielding diverse and customizable pluralistic environments.

600 7 Related Works

601 The human-like capabilities of LLMs have drawn significant attention from social science researchers, prompt-
602 ing extensive exploration at the intersection of AI and social sciences (Xu et al., 2024a). A key development
603 in this area is the shift from traditional Agent-Based Modeling to LLM-based agents, as explained by Ma
604 et al. (2024) through computational experiments. Numerous studies have since applied LLM-based agents
605 to diverse game scenarios, such as poker, Minecraft, and DOTA II, with more detailed summaries provided
606 by (Xu et al., 2024b; Hu et al., 2024b;a). Furthermore, Zhang et al. (2024c) have analyzed the core strategic
607 reasoning capabilities of these agents, distinguishing them from other reasoning approaches. While the previ-
608 ous reviews provide comprehensive overviews of related fields, our survey specifically focuses on social agents
609 equipped with beliefs, preferences, and reasoning capabilities within diverse game-theoretic scenarios.

610 8 Conclusion

611 We provide a comprehensive summary of existing research on LLM-based social agents in game-theoretic
612 scenarios from three perspectives: game framework, social agents, and evaluation protocol. This interdis-
613 ciplinary field covers a wide range of topics, including social sciences, economics, decision sciences, and theory
614 of mind. Current studies have primarily explored the more direct external behavioural patterns and internal
615 cognition of social agents. Therefore, future research should focus on developing theoretical frameworks
616 for cognitive representations within LLMs, conducting in-depth analyses of implicit and long-term game
617 behaviour patterns, and enhancing agents’ reasoning and planning capabilities in dynamic environments.

618 Broader Impact Statement

619 Developing agents with advanced social intelligence is one of the ultimate goals of artificial intelligence. On
620 one hand, such agents demonstrate enhanced collaboration, a deeper understanding of mental states, and
621 seamless integration into human society. On the other hand, negative social behaviors may also emerge, such
622 as deception, malicious competition, and verbal aggression, which conflict with the vision of a harmonious
623 human-AI coexistence.

Therefore, we carefully examine the potential negative impacts that social agents may have on human society, serving as a cautionary perspective for future social agent development. One major concern is *deception and manipulation*, where agents may bluff or mislead to achieve strategic goals. They may also engage in *malicious competition*, exploiting others to gain advantage, or exhibit *verbal and social aggression*, such as generating insults or polarizing language. Additionally, social agents can *amplify societal biases*, leading to discriminatory behaviors, and contribute to the *erosion of trust*, especially when users struggle to distinguish genuine human interactions from artificial ones. These agents may further *undermine human autonomy* by subtly steering decisions through persuasion, often without transparency. Due to their *scalability of harm*, even a single flawed agent can rapidly propagate misinformation or harmful behaviors across platforms. Moreover, the risk of *impersonation and infiltration* arises when agents mimic human users, potentially deceiving communities or individuals. These challenges highlight the critical need for careful design, value alignment, and robust supervision in the development and deployment of socially intelligent agents.

We now categorize the development and deployment of social agents into four stages: (1) Designing social agents, (2) Evaluating social agents, (3) Deploying social agents, and (4) Supervising social agents. Accordingly, we discuss the potential risks and feasible mitigation strategies for each stage. *Design Phase*: The underlying algorithms determine the agent’s behavioral preferences. Poorly designed algorithms may inadvertently lead to negative behaviors. To address this, researchers should enhance alignment algorithms, including safety alignment and moral alignment, to mitigate these risks at a fundamental level. Another promising approach is the design of behavioral plugins, where small models trained as plug-and-play behavior controllers can regulate agent actions dynamically. *Evaluation Phase*: Rigorous evaluation is crucial before deploying social agents in real-world applications. Agents exhibiting negative behaviors should be prevented from entering the deployment phase. One effective approach is to evaluate social agents across diverse game scenarios, allowing for a benchmarking framework that assesses their behavioral preferences under dynamic conditions. *Deployment Phase*: Direct large-scale deployment may lead to unforeseen negative consequences that were not observed in smaller-scale testing. Therefore, social agents should first be deployed in low-risk, small-scale environments, with a gradual expansion in scope and scale to monitor anomalies in real time. *Supervision Phase*: Effective oversight of social agents is essential. This can be achieved by designing automated monitoring systems that enable large-scale real-time surveillance. Behavioral analysis can be used to issue early warnings, assisting human supervisors in decision-making.

Additionally, it is important to note that most of the studies referenced in this paper utilize the GPT series as the large language model, which limits the generalizability of the experimental results. Differences in model architectures, training data, and alignment techniques can significantly impact the behavioral patterns exhibited by different models. Future research should explore a broader range of large language models, such as Claude, Gemini, Llama, and DeepSeek, to derive more comprehensive and reliable conclusions.

References

- Sahar Abdelnabi, Amr Gomaa, Sarath Sivaprasad, Lea Schönherr, and Mario Fritz. Llm-deliberation: Evaluating llms with interactive multi-agent negotiation games. *ArXiv preprint*, abs/2309.17234, 2023. URL <https://arxiv.org/abs/2309.17234>.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *ArXiv preprint*, abs/2303.08774, 2023. URL <https://arxiv.org/abs/2303.08774>.
- Gati Aher, RosaI. Arriaga, and Adam Tauman Kalai. Using large language models to simulate multiple humans and replicate human subject studies. In *International Conference on Machine Learning*, 2022. URL <https://api.semanticscholar.org/CorpusID:251719353>.
- Elif Akata, Lion Schulz, Julian Coda-Forno, Seong Joon Oh, Matthias Bethge, and Eric Schulz. Playing repeated games with large language models. *ArXiv preprint*, abs/2305.16867, 2023. URL <https://arxiv.org/abs/2305.16867>.
- Christopher Amato. An introduction to centralized training for decentralized execution in cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:2409.03052*, 2024.

- 673 Martin Atzmueller. Data mining on social interaction networks. Journal of Data Mining & Digital
674 Humanities, 2014, 2014.
- 675 Suma Bailis, Jane Friedhoff, and Feiyang Chen. Werewolf arena: A case study in llm evaluation via social
676 deduction. ArXiv preprint, abs/2407.13943, 2024. URL <https://arxiv.org/abs/2407.13943>.
- 677 Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff,
678 Jonathan Gray, Hengyuan Hu, Athul Paul Jacob, Mojtaba Komeili, Karthik Konath, Minae Kwon, Adam
679 Lerer, Mike Lewis, Alexander H. Miller, Sandra Mitts, Adithya Renduchintala, Stephen Roller, Dirk Rowe,
680 Weiyang Shi, Joe Spisak, Alexander Wei, David J. Wu, Hugh Zhang, and Markus Zijlstra. Human-level
681 play in the game of diplomacy by combining language models with strategic reasoning. Science, 378:1067
682 – 1074, 2022. URL <https://api.semanticscholar.org/CorpusID:253759631>.
- 683 Michael L Barnes and Robert J Sternberg. Social intelligence and decoding of nonverbal cues. Intelligence,
684 13(3):263–287, 1989.
- 685 Max H Bazerman, Jared R Curhan, Don A Moore, and Kathleen L Valley. Negotiation. Annual review of
686 psychology, 51(1):279–314, 2000.
- 687 Federico Bianchi, Patrick John Chia, Mert Yuksekgonul, Jacopo Tagliabue, Dan Jurafsky, and James Zou.
688 How well can llms negotiate? negotiationarena platform and analysis. ArXiv preprint, abs/2402.05863,
689 2024. URL <https://arxiv.org/abs/2402.05863>.
- 690 Matteo Bortoletto, Constantin Ruhdorfer, Lei Shi, and Andreas Bulling. Benchmarking mental state rep-
691 resentations in language models. ArXiv preprint, abs/2406.17513, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2406.17513)
692 [2406.17513](https://arxiv.org/abs/2406.17513).
- 693 Philip Brookins and Jason Matthew DeBacker. Playing games with gpt: What can we learn about a large
694 language model from canonical strategic games? Available at SSRN 4493398, 2023.
- 695 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter
696 Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early
697 experiments with gpt-4. ArXiv preprint, abs/2303.12712, 2023. URL [https://arxiv.org/abs/2303.](https://arxiv.org/abs/2303.12712)
698 [12712](https://arxiv.org/abs/2303.12712).
- 699 Colin F Camerer. Behavioral game theory: Experiments in strategic interaction. Princeton university press,
700 2011.
- 701 Jiangjie Chen, Siyu Yuan, Rong Ye, Bodhisattwa Prasad Majumder, and Kyle Richardson. Put your money
702 where your mouth is: Evaluating strategic planning and execution of llm agents in an auction arena. ArXiv
703 preprint, abs/2310.05746, 2023. URL <https://arxiv.org/abs/2310.05746>.
- 704 Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang,
705 Tinghui Zhu, et al. From persona to personalization: A survey on role-playing language agents. ArXiv
706 preprint, abs/2404.18231, 2024a. URL <https://arxiv.org/abs/2404.18231>.
- 707 Zehui Chen, Kuikun Liu, Qiuchen Wang, Jiangning Liu, Wenwei Zhang, Kai Chen, and Feng Zhao. Mind-
708 search: Mimicking human minds elicits deep ai searcher. ArXiv preprint, abs/2407.20183, 2024b. URL
709 <https://arxiv.org/abs/2407.20183>.
- 710 Richard Cheng, Abhinav Verma, Gabor Orosz, Swarat Chaudhuri, Yisong Yue, and Joel Burdick. Con-
711 trol regularization for reduced variance reinforcement learning. In International Conference on Machine
712 Learning, pp. 1141–1150. PMLR, 2019.
- 713 Yizhou Chi, Lingjun Mao, and Zineng Tang. Amongagents: Evaluating large language models in the
714 interactive text-based social deduction game. ArXiv preprint, abs/2407.16521, 2024. URL [https:](https://arxiv.org/abs/2407.16521)
715 [//arxiv.org/abs/2407.16521](https://arxiv.org/abs/2407.16521).
- 716 OpenCompass Contributors. Opencompass: A universal evaluation platform for foundation models. [https:](https://github.com/open-compass/opencompass)
717 [//github.com/open-compass/opencompass](https://github.com/open-compass/opencompass), 2023.

- Anthony Costarelli, Mat Allen, Roman Hauksson, Grace Sodunke, Suhas Hariharan, Carlson Cheng, Wenjie Li, and Arjun Yadav. Gamebench: Evaluating strategic reasoning abilities of llm agents. ArXiv preprint, abs/2406.06613, 2024. URL <https://arxiv.org/abs/2406.06613>. 718
719
720
- Daantje Derks, Arjan ER Bos, and Jasper Von Grumbkow. Emoticons and social interaction on the internet: the importance of social context. Computers in human behavior, 23(1):842–849, 2007. 721
722
- Jinhao Duan, Shiqi Wang, James Diffenderfer, Lichao Sun, Tianlong Chen, Bhavya Kailkhura, and Kaidi Xu. Reta: Recursively thinking ahead to improve the strategic reasoning of large language models. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pp. 2232–2246, 2024a. 723
724
725
726
- Jinhao Duan, Renming Zhang, James Diffenderfer, Bhavya Kailkhura, Lichao Sun, Elias Stengel-Eskin, Mohit Bansal, Tianlong Chen, and Kaidi Xu. Gtbench: Uncovering the strategic reasoning limitations of llms via game-theoretic evaluations. ArXiv preprint, abs/2402.12348, 2024b. URL <https://arxiv.org/abs/2402.12348>. 727
728
729
730
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. ArXiv preprint, abs/2407.21783, 2024. URL <https://arxiv.org/abs/2407.21783>. 731
732
733
- Caoyun Fan, Jindou Chen, Yaohui Jin, and Hao He. Can large language models serve as rational players in game theory? a systematic analysis. ArXiv preprint, abs/2312.05488, 2023. URL <https://arxiv.org/abs/2312.05488>. 734
735
736
- Xidong Feng, Yicheng Luo, Ziyang Wang, Hongrui Tang, Mengyue Yang, Kun Shao, David Mguni, Yali Du, and Jun Wang. Chessgpt: Bridging policy learning and language modeling. Advances in Neural Information Processing Systems, 36, 2024. 737
738
739
- Chaim Fershtman. The importance of the agenda in bargaining. Games and Economic Behavior, 2(3): 224–238, 1990. 740
741
- Nicol’o Fontana, Francesco Pierri, and Luca Maria Aiello. Nicer than humans: How do large language models behave in the prisoner’s dilemma? ArXiv preprint, abs/2406.13605, 2024. URL <https://arxiv.org/abs/2406.13605>. 742
743
744
- Martin E Ford and Marie S Tisak. A further search for social intelligence. Journal of Educational Psychology, 75(2):196, 1983. 745
746
- Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. Improving language model negotiation with self-play and in-context learning from ai feedback. ArXiv preprint, abs/2305.10142, 2023. URL <https://arxiv.org/abs/2305.10142>. 747
748
749
- Drew Fudenberg and Jean Tirole. Game theory. MIT press, 1991. 750
- Kanishk Gandhi, Dorsa Sadigh, and Noah D Goodman. Strategic reasoning with language models. ArXiv preprint, abs/2305.19165, 2023. URL <https://arxiv.org/abs/2305.19165>. 751
752
- Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. Understanding social reasoning in language models with language models. Advances in Neural Information Processing Systems, 36, 2024. 753
754
755
- Michael Georgeff, Barney Pell, Martha Pollack, Milind Tambe, and Michael Wooldridge. The belief-desire-intention model of agency. In Intelligent Agents V: Agents Theories, Architectures, and Languages: 5th International Workshop, ATAL’98 Paris, France, July 4–7, 1998 Proceedings 5, pp. 1–10. Springer, 1999. 756
757
758
- Fosca Giannotti, Mirco Nanni, Fabio Pinelli, and Dino Pedreschi. Trajectory pattern mining. In Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 330–339, 2007. 759
760
761

- 762 Alison Gopnik and Janet W Astington. Children’s understanding of representational change and its relation
763 to the understanding of false belief and the appearance-reality distinction. Child development, pp. 26–37,
764 1988.
- 765 Zhenyu Guan, Xiangyu Kong, Fangwei Zhong, and Yizhou Wang. Richelieu: Self-evolving llm-based agents
766 for ai diplomacy. ArXiv preprint, abs/2407.06813, 2024. URL <https://arxiv.org/abs/2407.06813>.
- 767 Fulin Guo. Gpt in game theory experiments. ArXiv preprint, abs/2305.05516, 2023. URL <https://arxiv.org/abs/2305.05516>.
- 769 Jiaxian Guo, Bo Yang, Paul Yoo, Bill Yuchen Lin, Yusuke Iwasawa, and Yutaka Matsuo. Suspicion-agent:
770 Playing imperfect information games with theory of mind aware gpt-4. ArXiv preprint, abs/2309.17277,
771 2023. URL <https://arxiv.org/abs/2309.17277>.
- 772 Shangmin Guo, Haoran Bu, Haochuan Wang, Yi Ren, Dianbo Sui, Yuming Shang, and Siting Lu. Economics
773 arena for large language models. ArXiv preprint, abs/2401.01735, 2024. URL <https://arxiv.org/abs/2401.01735>.
- 775 Akshat Gupta. Are chatgpt and gpt-4 good poker players?—a pre-flop analysis. ArXiv preprint,
776 abs/2308.12466, 2023. URL <https://arxiv.org/abs/2308.12466>.
- 777 Peter Hase, Thomas Hofweber, Xiang Zhou, Elias Stengel-Eskin, and Mohit Bansal. Fundamental problems
778 with model editing: How should rational belief revision work in llms? ArXiv preprint, abs/2406.19354,
779 2024. URL <https://arxiv.org/abs/2406.19354>.
- 780 He He, Jordan Boyd-Graber, Kevin Kwok, and Hal Daumé III. Opponent modeling in deep reinforcement
781 learning. In International conference on machine learning, pp. 1804–1813. PMLR, 2016.
- 782 Daniel A Herrmann and Benjamin A Levinstein. Standards for belief representations in llms. ArXiv preprint,
783 abs/2405.21030, 2024. URL <https://arxiv.org/abs/2405.21030>.
- 784 John J. Horton. Large language models as simulated economic agents: What can we learn from homo silicus?
785 SSRN Electronic Journal, 2023. URL <https://api.semanticscholar.org/CorpusID:255152420>.
- 786 Dirk Hovy and Diyi Yang. The importance of modeling social factors of language: Theory and practice. In
787 Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational
788 Linguistics: Human Language Technologies, pp. 588–602, Online, 2021. Association for Computational
789 Linguistics. doi: 10.18653/v1/2021.naacl-main.49. URL [https://aclanthology.org/2021.naacl-main.](https://aclanthology.org/2021.naacl-main.49)
790 [49](https://aclanthology.org/2021.naacl-main.49).
- 791 Chengpeng Hu, Yunlong Zhao, Ziqi Wang, Haocheng Du, and Jialin Liu. Games for artificial intelligence
792 research: A review and perspectives. IEEE Transactions on Artificial Intelligence, 2024a.
- 793 Sihao Hu, Tiansheng Huang, Fatih Ilhan, Selim Tekin, Gaowen Liu, Ramana Kompella, and Ling Liu.
794 A survey on large language model-based game agents. ArXiv preprint, abs/2404.02039, 2024b. URL
795 <https://arxiv.org/abs/2404.02039>.
- 796 Wenyue Hua, Ollie Liu, Lingyao Li, Alfonso Amayuelas, Julie Chen, Lucas Jiang, Mingyu Jin, Lizhou Fan,
797 Fei Sun, William Wang, et al. Game-theoretic llm: Agent workflow for negotiation games. arXiv preprint
798 arXiv:2411.05990, 2024a.
- 799 Yuncheng Hua, Lizhen Qu, and Gholamreza Haffari. Assistive large language model agents for socially-
800 aware negotiation dialogues. ArXiv preprint, abs/2402.01737, 2024b. URL <https://arxiv.org/abs/2402.01737>.
- 802 Chenghao Huang, Yanbo Cao, Yinlong Wen, Tao Zhou, and Yanru Zhang. Pokergpt: An end-to-end
803 lightweight solver for multi-player texas hold’em via large language model. ArXiv preprint, abs/2401.06781,
804 2024. URL <https://arxiv.org/abs/2401.06781>.
- 805 Thelma Hunt. The measurement of social intelligence. Journal of Applied Psychology, 12(3):317, 1928.

- Jingru Jia, Zehua Yuan, Junhao Pan, Paul E McNamara, and Deming Chen. Decision-making behavior evaluation framework for llms under uncertain context. arXiv preprint arXiv:2406.05972, 2024. 806 807
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. ArXiv preprint, abs/2310.06825, 2023. URL <https://arxiv.org/abs/2310.06825>. 808 809 810
- Xuanfa Jin, Ziyang Wang, Yali Du, Meng Fang, Haifeng Zhang, and Jun Wang. Learning to discuss strategically: A case study on one night ultimate werewolf. ArXiv preprint, abs/2405.19946, 2024. URL <https://arxiv.org/abs/2405.19946>. 811 812 813
- Chani Jung, Dongkwan Kim, Jiho Jin, Jiseon Kim, Yeon Seonwoo, Yejin Choi, Alice Oh, and Hyunwoo Kim. Perceptions to beliefs: Exploring precursory inferences for theory of mind in large language models. ArXiv preprint, abs/2407.06004, 2024. URL <https://arxiv.org/abs/2407.06004>. 814 815 816
- John H Kagel and Dan Levin. The winner’s curse and public information in common value auctions. American economic review, 76(5):894–920, 1986. 817 818
- Nora Kassner, Oyvind Tafjord, Ashish Sabharwal, Kyle Richardson, Hinrich Schuetze, and Peter Clark. Language models with rationality. ArXiv preprint, abs/2305.14250, 2023. URL <https://arxiv.org/abs/2305.14250>. 819 820 821
- Amirhossein Kazemnejad, Milad Aghajohari, Eva Portelance, Alessandro Sordani, Siva Reddy, Aaron Courville, and Nicolas Le Roux. Vineppo: Unlocking rl potential for llm reasoning through refined credit assignment. arXiv preprint arXiv:2410.01679, 2024. 822 823 824
- John F Kihlstrom and Nancy Cantor. Social intelligence. Handbook of intelligence, 2:359–379, 2000. 825
- Minsoo Kim, Jongyoon Kim, Jihyuk Kim, and Seung-won Hwang. QuBE: Question-based belief enhancement for agentic LLM reasoning. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pp. 21403–21423, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.1193. URL <https://aclanthology.org/2024.emnlp-main.1193/>. 826 827 828 829 830
- Henry Kissinger. Diplomacy. In Geopolitics, pp. 114–115. Routledge, 2014. 831
- Michal Kosinski. Theory of mind might have spontaneously emerged in large language models. ArXiv preprint, abs/2302.02083, 2023. URL <https://arxiv.org/abs/2302.02083>. 832 833
- Tiffany H Kung, Morgan Cheatham, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, et al. Performance of chatgpt on usmle: potential for ai-assisted medical education using large language models. PLoS digital health, 2(2): e0000198, 2023. 834 835 836 837
- Bolin Lai, Hongxin Zhang, Miao Liu, Aryan Pariani, Fiona Ryan, Wenqi Jia, Shirley Anugrah Hayati, James M Rehg, and Diyi Yang. Werewolf among us: A multimodal dataset for modeling persuasion behaviors in social deduction games. ArXiv preprint, abs/2212.08279, 2022. URL <https://arxiv.org/abs/2212.08279>. 838 839 840 841
- Yihuai Lan, Zhiqiang Hu, Lei Wang, Yang Wang, De-Yong Ye, Peilin Zhao, Ee-Peng Lim, Hui Xiong, and Hao Wang. Llm-based agent society investigation: Collaboration and confrontation in avalon gameplay. ArXiv preprint, abs/2310.14985, 2023. URL <https://arxiv.org/abs/2310.14985>. 842 843 844
- Yan Leng and Yuan Yuan. Do llm agents exhibit social behavior? ArXiv preprint, abs/2312.15198, 2023. URL <https://arxiv.org/abs/2312.15198>. 845 846
- Benjamin A Levinstein and Daniel A Herrmann. Still no lie detector for language models: Probing empirical and conceptual roadblocks. Philosophical Studies, pp. 1–27, 2024. 847 848

- 849 Huao Li, Yu Quan Chong, Simon Stepputtis, Joseph Campbell, Dana Hughes, Michael Lewis, and Ka-
850 tia Sycara. Theory of mind for multi-agent collaboration via large language models. ArXiv preprint,
851 abs/2310.10701, 2023a. URL <https://arxiv.org/abs/2310.10701>.
- 852 Jiatong Li, Rui Li, and Qi Liu. Beyond static datasets: A deep interaction approach to llm evaluation.
853 ArXiv preprint, abs/2309.04369, 2023b. URL <https://arxiv.org/abs/2309.04369>.
- 854 Minzhi Li, Weiyang Shi, Caleb Ziems, and Diyi Yang. Social intelligence data infrastructure: Structuring the
855 present and navigating the future. ArXiv preprint, abs/2403.14659, 2024a. URL <https://arxiv.org/abs/2403.14659>.
- 857 Yang Li, Shao Zhang, Jichen Sun, Yali Du, Ying Wen, Xinbing Wang, and Wei Pan. Cooperative open-
858 ended learning framework for zero-shot coordination. In International Conference on Machine Learning,
859 pp. 20470–20484. PMLR, 2023c.
- 860 Yuanchun Li, Hao Wen, Weijun Wang, Xiangyu Li, Yizhen Yuan, Guohong Liu, Jiacheng Liu, Wenxing Xu,
861 Xiang Wang, Yi Sun, et al. Personal llm agents: Insights and survey about the capability, efficiency and
862 security. ArXiv preprint, abs/2401.05459, 2024b. URL <https://arxiv.org/abs/2401.05459>.
- 863 Zongyuan Li, Yanan Ni, Runnan Qi, Lumin Jiang, Chang Lu, Xiaojie Xu, Xiangbei Liu, Pengfei Li, Yunzheng
864 Guo, Zhe Ma, et al. Llm-pysc2: Starcraft ii learning environment for large language models. arXiv preprint
865 [arXiv:2411.05348](https://arxiv.org/abs/2411.05348), 2024c.
- 866 Austen Liao, Nicholas Tomlin, and Dan Klein. Efficacy of language model self-play in non-zero-sum games.
867 ArXiv preprint, abs/2406.18872, 2024. URL <https://arxiv.org/abs/2406.18872>.
- 868 Jonathan Light, Min Cai, Sheng Shen, and Ziniu Hu. Avalonbench: Evaluating llms playing the game of
869 avalon. In NeurIPS 2023 Foundation Models for Decision Making Workshop, 2023a.
- 870 Jonathan Light, Min Cai, Sheng Shen, and Ziniu Hu. From text to tactic: Evaluating llms playing the game
871 of avalon. ArXiv preprint, abs/2310.05036, 2023b. URL <https://arxiv.org/abs/2310.05036>.
- 872 Jonathan Light, Min Cai, Weiqin Chen, Guanzhi Wang, Xiusi Chen, Wei Cheng, Yisong Yue, and Ziniu Hu.
873 Strategist: Learning strategic skills by llms via bi-level tree search. ArXiv preprint, abs/2408.10635, 2024.
874 URL <https://arxiv.org/abs/2408.10635>.
- 875 Yang Liu, Weixing Chen, Yongjie Bai, Jingzhou Luo, Xinshuai Song, Kaixuan Jiang, Zhida Li, Ganlong
876 Zhao, Junyi Lin, Guanbin Li, et al. Aligning cyber space with physical world: A comprehensive survey
877 on embodied ai. ArXiv preprint, abs/2407.06886, 2024a. URL <https://arxiv.org/abs/2407.06886>.
- 878 Ziyi Liu, Abhishek Anand, Pei Zhou, Jen-tse Huang, and Jieyu Zhao. Interintent: Investigating social intelli-
879 gence of llms via intention understanding in an interactive game context. ArXiv preprint, abs/2406.12203,
880 2024b. URL <https://arxiv.org/abs/2406.12203>.
- 881 Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao Ding, Gang Chen, and Haobo Wang. On llms-driven
882 synthetic data generation, curation, and evaluation: A survey. ArXiv preprint, abs/2406.15126, 2024.
883 URL <https://arxiv.org/abs/2406.15126>.
- 884 Nunzio Lorè and Babak Heydari. Strategic behavior of large language models and the role of game structure
885 versus contextual framing. Scientific Reports, 14(1):18490, 2024.
- 886 Ji Ma. Can machines think like humans? a behavioral evaluation of llm-agents in dictator games. arXiv
887 preprint [arXiv:2410.21359](https://arxiv.org/abs/2410.21359), 2024.
- 888 Qun Ma, Xiao Xue, Deyu Zhou, Xiangning Yu, Donghua Liu, Xuwen Zhang, Zihan Zhao, Yifan Shen, Peilin
889 Ji, Juanjuan Li, et al. Computational experiments meet large language model based agents: A survey and
890 perspective. ArXiv preprint, abs/2402.00262, 2024. URL <https://arxiv.org/abs/2402.00262>.

- Weiyu Ma, Qirui Mi, Xue Yan, Yuqiao Wu, Runji Lin, Haifeng Zhang, and Jun Wang. Large language models play starcraft ii: Benchmarks and a chain of summarization approach. ArXiv preprint, abs/2312.11865, 2023. URL <https://arxiv.org/abs/2312.11865>.
- Bahar Mahmud, Guan Hong, and Bernard Fong. A study of human–ai symbiosis for creative work: Recent developments and future directions in deep learning. ACM Transactions on Multimedia Computing, Communications and Applications, 20(2):1–21, 2023.
- Shaoguang Mao, Yuzhe Cai, Yan Xia, Wenshan Wu, Xun Wang, Fengyi Wang, Tao Ge, and Furu Wei. Alympics: Language agents meet game theory. ArXiv preprint, abs/2311.03220, 2023. URL <https://arxiv.org/abs/2311.03220>.
- Leena Mathur, Paul Pu Liang, and Louis-Philippe Morency. Advancing social intelligence in ai agents: Technical challenges and open questions. ArXiv preprint, abs/2404.11023, 2024. URL <https://arxiv.org/abs/2404.11023>.
- Qiaozhu Mei, Yutong Xie, Walter Yuan, and Matthew O Jackson. A turing test of whether ai chatbots are behaviorally similar to humans. Proceedings of the National Academy of Sciences, 121(9):e2313925121, 2024.
- Juanjuan Meng. Ai emerges as the frontier in behavioral science. Proceedings of the National Academy of Sciences of the United States of America, 121 10:e2401336121, 2024. URL <https://api.semanticscholar.org/CorpusID:268029379>.
- Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. Science, 356(6337):508–513, 2017.
- Rémi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland, Zhao-han Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Andrea Michi, et al. Nash learning from human feedback. ArXiv preprint, abs/2312.00886, 2023. URL <https://arxiv.org/abs/2312.00886>.
- Aran Nayebi. Barriers and pathways to human-ai alignment: A game-theoretic approach. arXiv preprint arXiv:2502.05934, 2025.
- Sean Noh and Ho-Chun Herbert Chang. Llms with personalities in multi-issue negotiation games. ArXiv preprint, abs/2405.05248, 2024. URL <https://arxiv.org/abs/2405.05248>.
- Maayan Orner, Oleg Maksimov, Akiva Kleinerman, Charles Ortiz, and Sarit Kraus. Explaining decisions of agents in mixed-motive games. ArXiv preprint, abs/2407.15255, 2024. URL <https://arxiv.org/abs/2407.15255>.
- Guillermo Owen. Game theory. Emerald Group Publishing, 2013.
- Steve Phelps and Yvan I. Russell. The machine psychology of cooperation: Can gpt models operationalise prompts for altruism, cooperation, competitiveness and selfishness in economic games? ArXiv preprint, 2023. URL <https://api.semanticscholar.org/CorpusID:258685424>.
- Giorgio Piatti, Zhijing Jin, Max Kleiman-Weiner, Bernhard Schölkopf, Mrinmaya Sachan, and Rada Mihalcea. Cooperate or collapse: Emergence of sustainability behaviors in a society of llm agents. ArXiv preprint, abs/2404.16698, 2024. URL <https://arxiv.org/abs/2404.16698>.
- David Premack and Guy Woodruff. Does the chimpanzee have a theory of mind? Behavioral and brain sciences, 1(4):515–526, 1978.
- Siyuan Qi, Shuo Chen, Yexin Li, Xiangyu Kong, Junqi Wang, Bangcheng Yang, Pring Wong, Yifan Zhong, Xiaoyuan Zhang, Zhaowei Zhang, Nian Liu, Wei Wang, Yaodong Yang, and Song-Chun Zhu. Civrealm: A learning and reasoning odyssey in civilization for decision-making agents. ArXiv preprint, abs/2401.10568, 2024. URL <https://arxiv.org/abs/2401.10568>.

- 935 Huida Qiu, Yan Liu, Niranjan A Subrahmanya, and Weichang Li. Granger causality for time-series anomaly
936 detection. In 2012 IEEE 12th international conference on data mining, pp. 1074–1079. IEEE, 2012.
- 937 Anatol Rapoport and Albert M Chammah. Prisoner’s dilemma: A study in conflict and cooperation, volume
938 165. University of Michigan press, 1965.
- 939 Reza Rawassizadeh, Elaheh Momeni, Chelsea Dobbins, Joobin Gharibshah, and Michael Pazzani. Scalable
940 daily human behavioral pattern mining from multivariate temporal data. IEEE Transactions on Knowledge
941 and Data Engineering, 28(11):3098–3112, 2016.
- 942 Siyue Ren, Zhiyao Cui, Ruiqi Song, Zhen Wang, and Shuyue Hu. Emergence of social norms in generative
943 agent societies: principles and architecture. In Proceedings of the 33rd International Joint Conference on
944 Artificial Intelligence (IJCAI), 2024.
- 945 Jillian Ross, Yoon Kim, and Andrew W Lo. Llm economicus? mapping the behavioral biases of llms via
946 utility theory. ArXiv preprint, abs/2408.02784, 2024. URL <https://arxiv.org/abs/2408.02784>.
- 947 Nino Scherrer, Claudia Shi, Amir Feder, and David Blei. Evaluating the moral beliefs encoded in llms.
948 Advances in Neural Information Processing Systems, 36, 2024.
- 949 Stefan F Schouten, Peter Bloem, Ilia Markov, and Piek Vossen. Truth-value judgment in language models:
950 belief directions are context sensitive. ArXiv preprint, abs/2404.18865, 2024. URL [https://arxiv.org/
951 abs/2404.18865](https://arxiv.org/abs/2404.18865).
- 952 Melanie Sclar, Sachin Kumar, Peter West, Alane Suhr, Yejin Choi, and Yulia Tsvetkov. Minding lan-
953 guage models’(lack of) theory of mind: A plug-and-play multi-character belief tracker. ArXiv preprint,
954 abs/2306.00924, 2023. URL <https://arxiv.org/abs/2306.00924>.
- 955 Xiao Shao, Weifu Jiang, Fei Zuo, and Mengqing Liu. Swarmbrain: Embodied agent for real-time strategy
956 game starcraft ii via large language models. ArXiv preprint, abs/2401.17749, 2024. URL [https://arxiv.
957 org/abs/2401.17749](https://arxiv.org/abs/2401.17749).
- 958 Eilam Shapira, Omer Madmon, Itamar Reinman, Samuel Joseph Amouyal, Roi Reichart, and Moshe Ten-
959 nenholtz. Glee: A unified framework and benchmark for language-based economic environments. arXiv
960 preprint arXiv:2410.05254, 2024.
- 961 Zijing Shi, Meng Fang, Shunfeng Zheng, Shilong Deng, Ling Chen, and Yali Du. Cooperation on the fly:
962 Exploring language agents for ad hoc teamwork in the avalon game. ArXiv preprint, abs/2312.17515,
963 2023. URL <https://arxiv.org/abs/2312.17515>.
- 964 Hisaichi Shibata, Soichiro Miki, and Yuta Nakamura. Playing the werewolf game with artificial intelligence
965 for language understanding. arXiv preprint arXiv:2302.10646, 2023.
- 966 Theodore R Summers, Shunyu Yao, Karthik Narasimhan, and Thomas L Griffiths. Cognitive architectures for
967 language agents. ArXiv preprint, abs/2309.02427, 2023. URL <https://arxiv.org/abs/2309.02427>.
- 968 Lawrence E Susskind. Scorable games: A better way to teach negotiation. Negot. J., 1:205, 1985.
- 969 Reiji Suzuki and Takaya Arita. An evolutionary model of personality traits related to cooperative behavior
970 using a large language model. Scientific Reports, 14, 2023. URL [https://api.semanticscholar.org/
971 CorpusID:263830498](https://api.semanticscholar.org/CorpusID:263830498).
- 972 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Sori-
973 cut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal
974 models. ArXiv preprint, abs/2312.11805, 2023. URL <https://arxiv.org/abs/2312.11805>.
- 975 Jen tse Huang, Eric Li, Man Ho Lam, Tian Liang, Wenxuan Wang, Youliang Yuan, Wenxiang Jiao, Xing
976 Wang, Zhaopeng Tu, and Michael R. Lyu. How far are we on the decision-making of llms? evaluating
977 llms’ gaming ability in multi-agent environments. ArXiv preprint, abs/2403.11807, 2024. URL [https://
978 //arxiv.org/abs/2403.11807](https://arxiv.org/abs/2403.11807).

- Wiebe Van Der Hoek, Wojciech Jamroga, and Michael Wooldridge. A logic for strategic reasoning. In Proceedings of the fourth international joint conference on Autonomous agents and multiagent systems, pp. 157–164, 2005. 979–981
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. Frontiers of Computer Science, 18(6):186345, 2024a. 982–984
- Shenzhi Wang, Chang Liu, Zilong Zheng, Siyuan Qi, Shuo Chen, Qisen Yang, Andrew Zhao, Chaofei Wang, Shiji Song, and Gao Huang. Avalon’s game of thoughts: Battle against deception through recursive contemplation. ArXiv preprint, abs/2310.01320, 2023. URL <https://arxiv.org/abs/2310.01320>. 985–987
- Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Opendevin: An open platform for ai software developers as generalist agents. ArXiv preprint, abs/2407.16741, 2024b. URL <https://arxiv.org/abs/2407.16741>. 988–990
- Zhen Wang, Ruiqi Song, Chen Shen, Shiya Yin, Zhao Song, Balaraju Battu, Lei Shi, Danyang Jia, Talal Rahwan, and Shuyue Hu. Large language models overcome the machine penalty when acting fairly but not when acting selfishly or altruistically. arXiv preprint arXiv:2410.03724, 2024c. 991–993
- Donald Arthur Waterman. Generalization learning techniques for automating the learning of heuristics. Artificial Intelligence, 1(1-2):121–170, 1970. 994–995
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022. 996–998
- Dekun Wu, Haochen Shi, Zhiyuan Sun, and Bang Liu. Deciphering digital detectives: Understanding llm behaviors and capabilities in multi-agent mystery games. In Findings of the Association for Computational Linguistics ACL 2024, pp. 8225–8291, 2024a. 999–1001
- Shuang Wu, Liwen Zhu, Tao Yang, Shiwei Xu, Qiang Fu, Yang Wei, and Haobo Fu. Enhance reasoning for large language models in the game werewolf. ArXiv preprint, abs/2402.02330, 2024b. URL <https://arxiv.org/abs/2402.02330>. 1002–1004
- Tian Xia, Zhiwei He, Tong Ren, Yibo Miao, Zhuosheng Zhang, Yang Yang, and Rui Wang. Measuring bargaining abilities of llms: A benchmark and a buyer-enhancement method. ArXiv preprint, abs/2402.15813, 2024. URL <https://arxiv.org/abs/2402.15813>. 1005–1007
- Lin Xu, Zhiyuan Hu, Daquan Zhou, Hongyu Ren, Zhen Dong, Kurt Keutzer, See-Kiong Ng, and Jiashi Feng. Magic: Investigation of large language model powered multi-agent in cognition, adaptability, rationality and collaboration. In ICLR 2024 Workshop on Large Language Model (LLM) Agents, 2023a. 1008–1010
- Rongwu Xu, Brian S Lin, Shujian Yang, Tianqi Zhang, Weiyang Shi, Tianwei Zhang, Zhixuan Fang, Wei Xu, and Han Qiu. The earth is flat because...: Investigating llms’ belief towards misinformation via persuasive conversation. ArXiv preprint, abs/2312.09085, 2023b. URL <https://arxiv.org/abs/2312.09085>. 1011–1013
- Ruoxi Xu, Yingfei Sun, Mengjie Ren, Shiguang Guo, Ruotong Pan, Hongyu Lin, Le Sun, and Xianpei Han. Ai for social science and social science of ai: A survey. Information Processing & Management, 61(3): 103665, 2024a. 1014–1016
- Xinrun Xu, Yuxin Wang, Chaoyi Xu, Ziluo Ding, Jiechuan Jiang, Zhiming Ding, and Börje F Karlsson. A survey on game playing agents and large models: Methods, applications, and challenges. ArXiv preprint, abs/2403.10249, 2024b. URL <https://arxiv.org/abs/2403.10249>. 1017–1019
- Yuzhuang Xu, Shuo Wang, Peng Li, Fuwen Luo, Xiaolong Wang, Weidong Liu, and Yang Liu. Exploring large language models for communication games: An empirical study on werewolf. ArXiv preprint, abs/2309.04658, 2023c. URL <https://arxiv.org/abs/2309.04658>. 1020–1022

- 1023 Zelai Xu, Chao Yu, Fei Fang, Yu Wang, and Yi Wu. Language agents with reinforcement learning for
1024 strategic play in the werewolf game. ArXiv preprint, abs/2310.18940, 2023d. URL <https://arxiv.org/abs/2310.18940>.
1025
- 1026 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li,
1027 Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. ArXiv preprint, abs/2407.10671, 2024a. URL
1028 <https://arxiv.org/abs/2407.10671>.
- 1029 Ruihan Yang, Jiangjie Chen, Yikai Zhang, Siyu Yuan, Aili Chen, Kyle Richardson, Yanghua Xiao, and
1030 Deqing Yang. Selfgoal: Your language agents already know how to achieve high-level goals. ArXiv
1031 preprint, abs/2406.04784, 2024b. URL <https://arxiv.org/abs/2406.04784>.
- 1032 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
1033 Tree of thoughts: Deliberate problem solving with large language models. Advances in Neural Information
1034 Processing Systems, 36, 2024.
- 1035 Yauwai Yim, Chunkit Chan, Tianyu Shi, Zheye Deng, Wei Fan, Tianshi Zheng, and Yangqiu Song. Evaluating
1036 and enhancing llms agent based on theory of mind in guandan: A multi-player cooperative game under
1037 imperfect information. ArXiv preprint, abs/2408.02559, 2024. URL [https://arxiv.org/abs/2408.](https://arxiv.org/abs/2408.02559)
1038 [02559](https://arxiv.org/abs/2408.02559).
- 1039 Haolan Zhan, Yufei Wang, Tao Feng, Yuncheng Hua, Suraj Sharma, Zhuang Li, Lizhen Qu, Zhaleh Semnani
1040 Azad, Ingrid Zukerman, and Gholamreza Haffari. Let’s negotiate! a survey of negotiation dialogue systems.
1041 ArXiv preprint, abs/2402.01097, 2024. URL <https://arxiv.org/abs/2402.01097>.
- 1042 Wenqi Zhang, Ke Tang, Hai Wu, Mengna Wang, Yongliang Shen, Guiyang Hou, Zeqi Tan, Peng Li, Yueting
1043 Zhuang, and Weiming Lu. Agent-pro: Learning to evolve via policy-level reflection and optimization.
1044 arXiv preprint arXiv:2402.17574, 2024a.
- 1045 Yadong Zhang, Shaoguang Mao, Tao Ge, Xun Wang, Adrian de Wynter, Yan Xia, Wenshan Wu, Ting Song,
1046 Man Lan, and Furu Wei. Llm as a mastermind: A survey of strategic reasoning with large language
1047 models. ArXiv preprint, abs/2404.01230, 2024b. URL <https://arxiv.org/abs/2404.01230>.
- 1048 Yadong Zhang, Shaoguang Mao, Tao Ge, Xun Wang, Adrian de Wynter, Yan Xia, Wenshan Wu, Ting Song,
1049 Man Lan, and Furu Wei. Llm as a mastermind: A survey of strategic reasoning with large language
1050 models. ArXiv preprint, abs/2404.01230, 2024c. URL <https://arxiv.org/abs/2404.01230>.
- 1051 Yadong Zhang, Shaoguang Mao, Tao Ge, Xun Wang, Yan Xia, Man Lan, and Furu Wei. K-level reasoning
1052 with large language models. ArXiv preprint, abs/2402.01521, 2024d. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2402.01521)
1053 [2402.01521](https://arxiv.org/abs/2402.01521).
- 1054 Yadong Zhang, Shaoguang Mao, Wenshan Wu, Yan Xia, Tao Ge, Man Lan, and Furu Wei. Enhancing
1055 language model rationality with bi-directional deliberation reasoning. ArXiv preprint, abs/2407.06112,
1056 2024e. URL <https://arxiv.org/abs/2407.06112>.
- 1057 Qinlin Zhao, Jindong Wang, Yixuan Zhang, Yiqiao Jin, Kaijie Zhu, Hao Chen, and Xing Xie. Competeai:
1058 Understanding the competition dynamics in large language model-based agents. In Forty-first International
1059 Conference on Machine Learning, 2023. URL <https://api.semanticscholar.org/CorpusID:270357283>.
- 1060 Qinglin Zhu, Runcong Zhao, Jinhua Du, Lin Gui, and Yulan He. Player*: Enhancing llm-based multi-agent
1061 communication and interaction in murder mystery games. ArXiv preprint, abs/2404.17662, 2024a. URL
1062 <https://arxiv.org/abs/2404.17662>.
- 1063 Wentao Zhu, Zhining Zhang, and Yizhou Wang. Language models represent beliefs of self and others. ArXiv
1064 preprint, abs/2402.18496, 2024b. URL <https://arxiv.org/abs/2402.18496>.
- 1065 Richard Zhuang, Akshat Gupta, Richard Yang, Aniket Rahane, Zhengyu Li, and Gopala Anumanchipalli.
1066 Pokerbench: Training large language models to become professional poker players. arXiv preprint
1067 arXiv:2501.08328, 2025.