
GTBENCH: Uncovering the Strategic Reasoning Limitations of LLMs via Game-Theoretic Evaluations

Jinhao Duan^{1*} Renming Zhang^{2*} James Diffenderfer³ Bhavya Kailkhura³
Lichao Sun⁴ Elias Stengel-Eskin⁵ Mohit Bansal⁵ Tianlong Chen^{5,6,7†} Kaidi Xu^{1†}
¹Drexel University ²Boston University ³LLNL ⁴Lehigh University
⁵UNC Chapel Hill ⁶MIT ⁷Harvard University

Abstract

As Large Language Models (LLMs) are integrated into critical real-world applications, their strategic and logical reasoning abilities are increasingly crucial. This paper evaluates LLMs’ reasoning abilities in competitive environments through game-theoretic tasks, e.g., board and card games that require pure logic and strategic reasoning to compete with opponents. We first propose GTBENCH, a language-driven environment composing 10 widely-recognized tasks, across a comprehensive game taxonomy: complete versus incomplete information, dynamic versus static, and probabilistic versus deterministic scenarios. Then, we ❶ Characterize the game-theoretic reasoning of LLMs; and ❷ Perform LLM-vs.-LLM competitions as reasoning evaluation. We observe that ❶ LLMs have distinct behaviors regarding various gaming scenarios; for example, LLMs fail in complete and deterministic games yet they are competitive in probabilistic gaming scenarios; ❷ Most open-source LLMs, e.g., CodeLlama-34b-Instruct and Llama-2-70b-chat, are less competitive than commercial LLMs, e.g., GPT-4, in complex games, yet the recently released Llama-3-70b-Instruct makes up for this shortcoming. In addition, code-pretraining greatly benefits strategic reasoning, while advanced reasoning methods such as Chain-of-Thought (CoT) and Tree-of-Thought (ToT) do not always help. We further characterize the game-theoretic properties of LLMs, such as equilibrium and Pareto Efficiency in repeated games. Detailed error profiles are provided for a better understanding of LLMs’ behavior. We hope our research provides standardized protocols and serves as a foundation to spur further explorations in the strategic reasoning of LLMs.

1 Introduction

Large Language Models (LLMs) are increasingly being integrated into critical real-world applications, such as cybersecurity (Ameri et al., 2021; Aghaei et al., 2022), decision science (Jiang et al., 2023b), and finance (Wu et al., 2023). These areas involve advanced strategic thinking and logical reasoning skills, including the ability to foresee possible dangers and weaknesses (Yao et al., 2024b; Duan et al., 2024a), systematically examine difficulties, and make informed decisions based on provided evidence. However, evaluation environments that thoroughly assess these situations are not sufficiently explored.

There has been an emerging trend where LLMs are evaluated in various interactive role-playing environments, including collaborative environments such as CAMEL (Li et al., 2023), ReConcile (Chen et al., 2023), and competition environments such as Diplomacy (Bakhtin et al., 2022), Werewolf (Xu et al., 2023a), Avalon (Light et al., 2023; Stepputtis et al., 2023), multi-agent debate (Liang et al., 2023; Du et al., 2023; Chan et al., 2023; Xiong et al., 2023), board and card games (Duan et al.,

* Equal contribution.

† Correspondence to: Tianlong Chen tianlong@mit.edu, Kaidi Xu kx46@drexel.edu

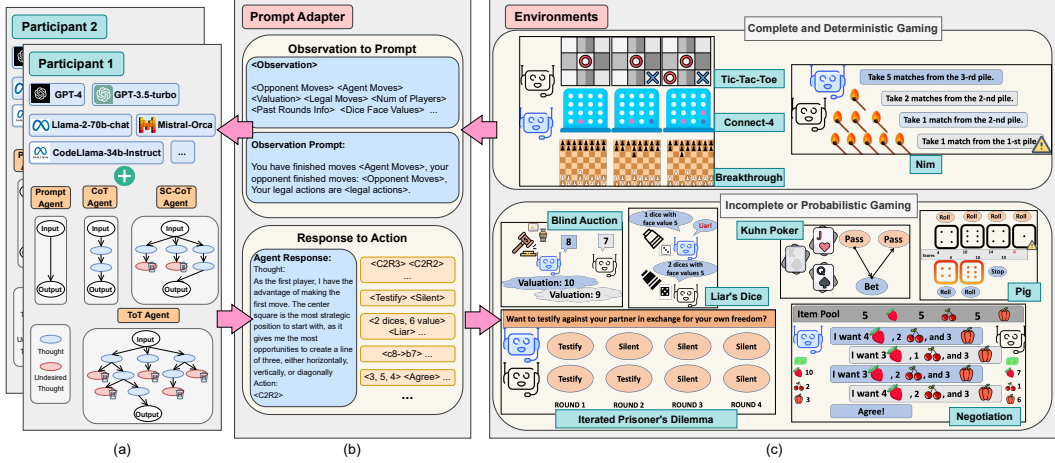


Figure 1: The overall schematic of GTBENCH. There are three main components from *right to left*: **Environments** (c) for game hosting, observation providing, and action execution; **Prompt Adapter** (b) for converting observation to prompt and extracting actions from participants’ generations; **Participants** (a) for reasoning and action generation.

2024b). By engaging LLMs in simulated scenarios, role-playing-based environments offer useful potential for analyzing the cognitive reasoning abilities of LLMs. However, the extensive background and intricate details involved in role-play-based games dilute the pureness of logic and strategic reasoning that is typically found in game-theoretic tasks. Additionally, the evaluation is primarily verbal as it hinges on spoken or written exchanges between the LLMs. This could mask instances where LLMs might *lack concrete reasoning abilities but navigate the scenario effectively through the proficient use of language*.

Why are game-theoretic tasks unique and necessary for LLM reasoning evaluation? Game-theoretic tasks are typically conceptualized based on prevalent trade-offs and dilemmas manifesting in real-life scenarios and are designed to be easy to understand yet require difficult skills to be mastered. In contrast to the rich narrative contexts afforded in verbal- or role-playing-based games, e.g., Werewolf (Xu et al., 2023a) and Avalon (Light et al., 2023), the reality of game-theoretic games such as Chess and Go involve: ❶ pure logic and strategic reasoning without the added complexity of backgrounds or character roles; ❷ embracing rigorous rules with well-defined action/state space, which allow for an in-depth examination of the strategic reasoning of LLMs.

Hence, in order to spur more research in the LLM Game-Theoretic evaluation domain, we propose GTBENCH, an environment consisting of 10 widely recognized game-theoretic tasks, across a comprehensive taxonomy of games, e.g., complete- (Tic-Tac-Toe, Connect-4, Breakthrough) versus incomplete-information (Kuhn Poker, Liar’s Dice) gaming, deterministic (Nim) versus probabilistic (Negotiation, Pig) gaming, static versus dynamic (Iterated Prisoner’s Dilemma, Blind Auction) gaming. These environments require a variety of abilities including board strategy, collaboration, auction, and bidding. There are two key issues investigated in this paper:

Characterizing Strategic Reasoning of LLMs: *How LLMs will perform when facing various game-theoretic scenarios? How do they perform compared to conventional solvers? How do essential factors, e.g., pertaining, parameter sizes, and reasoning methods, affect strategic reasoning?*

LLM-vs.-LLM Competitions as New Reasoning Evaluation: *A new automated and adaptive benchmark that can be effective in evaluating reasoning errors even for future LLMs.*

To address these crucial problems, we conduct experiments over two configurations: (a) **LLM-vs-Conventional** where conventional solvers such as optimization- or search-based solvers, e.g., Monte-Carlo Tree Search (MCTS) (Chaslot et al., 2008), are taken as the opponent of LLMs; (b) **LLM-vs.-LLM** where two LLMs compete directly to reveal the reasoning limitations in an automated manner. We find that: ❶ LLMs almost always fail when playing against simple MCTS opponents in complete and deterministic gaming scenarios (Section 4.1), while ❷ LLMs remain competitive in incomplete and probabilistic scenarios (Section 4.2); ❸ Code-pretraining benefits game-theoretic

reasoning, e.g., CodeLlama-34b-Instruct (Roziere et al., 2023) achieves comparable results as GPT-3.5-turbo, and significantly outperforms Llama-2-70b-chat (Touvron et al., 2023) (Section 4.3); ④ Advanced reasoning methods, such as Chain-of-Thought (CoT) (Wei et al., 2022), Self-Consistent CoT (SC-CoT) (Wang et al., 2022b), Tree-of-Thought (ToT) (Yao et al., 2024a) are not always helpful; ⑤ Most open-source LLMs are less competitive than commercial LLMs in games with complex rules and large action/state space, while the recently released Llama-3-70b-Instruct (Meta, 2024) makes up for this shortcoming. The interfaces of GTBENCH leaderboard can be found in Appendix A11. Our contributions can be summarized as the following:

- **LLM Game-Theoretic Evaluation (GTBENCH):** An LLM environment supporting 10 well-recognized tasks across comprehensive game-theoretic taxonomy, is presented to spur future work for the community. The code and leaderboard will be public and continuously updated for future reasoning agents and LLMs.
- **Essential Factors for the Strategic Reasoning of LLMs:** We investigate how essential factors, e.g., parameter size, code-pretraining, and reasoning methods, affect strategic reasoning. A detailed error profile is provided for a better understanding of LLMs’ behaviors.
- **Characteriz the Game-Theoretic Properties of LLMs:** We characterize distinct LLM behaviors when facing different game-theoretic scenarios, such as LLMs fail in complete-information and deterministic gaming yet remain competitive in probabilistic gaming. We further study the equilibrium and Pareto efficiency during the gameplay.

2 Background and Problem Definition

2.1 Background and Related Work

LLM-as-Agent Evaluation. Several studies have been conducted to measure the effectiveness of LLMs as agents in recent years. Hausknecht et al. (2020) carried out an extensive study to evaluate the performance of LLMs in interactive fiction games. Zhu et al. (2023) provides a valuable dataset for finetuning LLMs to improve usefulness in the strategic game Dungeons & Dragons. GRUE (Ramamurthy et al., 2023) uses reinforcement learning-based metrics to benchmark the performance of generation tasks in six different languages. Gandhi et al. (2023) test the use of LLMs as a broker with human contestants in the negotiation game “Deal or No Deal”. A few studies have explored the use of text-based games as a means of facilitating learning in such environments. ALFWorld (Shridhar et al., 2020) introduced a novel virtual environment that allows agents to acquire learning in a text-based environment while executing in a visual environment. The environment was developed in conjunction with Building Understanding in Text world via Language for Embodied Reasoning (BUTLER) agent, which can acquire abstract text knowledge in the text world. Similarly, TextWorld (Côté et al., 2019) is introduced as an environment that enables RL agents to play text games. Wang et al. (2022a) proposed ScienceWorld, a benchmark used for evaluating agents’ reasoning ability, and their findings showed that transformer-based models are not effective at reasoning in novel contexts. MTBench (Zheng et al., 2024) introduces LLM-as-a-Judge where GPT-4 (Achiam et al., 2023) is utilized as a judge to evaluate the quality of LLM generations. It indicates that GPT-4 shares close criteria as humans. There have been works evaluating LLMs in solving real-world tasks, such as graph reasoning (Besta et al., 2023), WebShop (Yao et al., 2022), AgentBench (Liu et al., 2023) for pragmatic missions, MINT (Wang et al., 2023b) for tool utilization.

Multiple LLMs-as-Agents in Gaming. A key research area is the competition and collaboration between LLMs. Many studies examine LLMs’ strategic reasoning and performance, using evaluation frameworks to assess multiple LLM agents in individual games, such as: Social deduction or deception games (Xu et al., 2023a,b; O’Gara, 2023; Light et al., 2023), diplomacy games (Mukobi et al., 2023; , FAIR), negotiation games (Abdelnabi et al., 2023; Davidson et al., 2023), coordination and cooperation games (Akata et al., 2023), and Minecraft (Gong et al., 2023; Wang et al., 2023a; Fan et al., 2022). These works not only provide evaluation frameworks for games and demonstrate the flexibility of LLMs to a variety of gaming tasks but some provide meaningful datasets for fine-tuning, policies for reinforcement learning to produce better strategies, or evaluate the strategic reasoning of LLMs. However, many of these standalone works quantify either individual or a subset of desirable strategic reasoning capabilities of LLMs, such as negotiation, deception, or coordination. Further, they often evaluate these capabilities for LLMs using one or two games which may produce less robust assurances of LLM abilities.

Table 1: Game environments explored in GTBENCH.

Game	Taxonomy of Games						Preferred Ability				
	Zero-Sum	First-player Advantage	▲ Complete ● Incomplete	▲ Dynamic ● Static	▲ Probabilistic ● Deterministic		Board Strategy	Bids	Collaboration	Bluff	Math
Tic-Tac-Toe	✓	✓	▲	●	●		✓	×	×	×	×
Connect-4	✓	✓	▲	●	●		✓	×	×	×	×
Kuhn Poker	✓	✓	●	●	▲		×	×	×	✓	✓
Breakthrough	✓	× [†]	▲	●	●		✓	×	×	×	×
Liar's Dice	✓	×	●	●	▲		×	✓	×	✓	✓
Blind Auction	×	×	●	▲	▲		×	✓	×	×	✓
Negotiation	×	×	●	●	▲		×	×	✓	✓	✓
Nim	✓	✓	▲	●	●		×	×	×	×	✓
Pig	×	×	▲	●	▲		×	×	×	×	×
Iterated Prisoner's Dilemma	×	×	▲	▲	●		×	×	✓ [‡]	×	✓

[†]: Breakthrough has a slight first-player advantage which is not as significant as others.

[‡]: The iterated version of Prisoner's Dilemma allows participants access to the actions made by their opponents in the past rounds, achieving implicit collaboration.

^{††}: Inapplicable due to complex combination and dynamic environment.

We make an additional crucial contribution in this line of work by measuring strategic reasoning capabilities with games that are not found in the existing unified benchmark suites (Zhang et al., 2024), such as clembench (Chalamalasetti et al., 2023) focusing on conversational agents over non-zero-sum games and LMRL-Gym (Abdulhai et al., 2023) on verbal reinforcement learning tasks. (Chen et al., 2024) and Duan et al. (2024b) also proposes multi-agent strategic reasoning evaluation. However, they overlooked the analysis of LLM behaviors in response to different game-theoretic scenarios and their associated properties. Differently, GTBENCH seeks to provide a unified suite of games that are carefully curated to (1) evaluate a comprehensive collection of strategic reasoning abilities for a given agent and (2) enable competition-based scenarios (i.e., LLM agent-1 vs LLM agent-2) allowing for competition-based comparisons of strategic reasoning capabilities by LLM-based agents.

2.2 Problem Definition

Notation: Gameplay. We formulate the gameplay as a Markov Decision Process $(\mathcal{S}, \mathcal{A}, \mathcal{M}, \mathcal{O})$ under a given game environment, among the alternating interaction of two participants. This process composes of an infinite state space \mathcal{S} , an infinite action space \mathcal{A} , the participants $\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2\}$, and an observation space \mathcal{O} . Considering the decision of \mathcal{M}_i ($i = 1, 2$) at the t -th step of the process, we denote by $s_t \in \mathcal{S}$ the state that \mathcal{M}_i are placed and $o_t \in \mathcal{O}$ the observation that \mathcal{M}_i are observing. We assume \mathcal{M}_i follows policy $\pi_{\theta_i}(a_t|s_t, o_t)$ for state transition $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$, where $a_t \in \mathcal{A}$ is the action sampled by π_{θ_i} under conditions s_t and o_t . θ_i is determined by the implementation by \mathcal{M}_i , e.g., optimization-based solver, LLM-driven agents, which will be discussed in Section 3.2 in detail. In this way, the two-participate gameplay can be represented as $(s_0, a_0, s_1, a_1, s_2, \dots, s_n)$, where s_0 is the initial state and s_n is a terminal state, i.e., end of the game. The progress is driven by the alternating execution of actions sampled by participants. Please refer to Section 3.1 and Appendix A2 for all the supported games with the corresponding actions and observations.

Evaluation Metric: Normalized Relative Advantage. We introduce **Normalized Relative Advantage (NRA)**, denoted $NRA(\mathcal{M}_i, \mathcal{M}_o, f_s)$, to measure to relative advantage of \mathcal{M}_i when competing against \mathcal{M}_o , under the score calculation f_s :

$$NRA(\mathcal{M}_i, \mathcal{M}_o, f_s) = \frac{\sum_m f_s(\mathcal{M}_i, m) - \sum_m f_s(\mathcal{M}_o, m)}{\sum_m f_s(\mathcal{M}_i, m) + \sum_m f_s(\mathcal{M}_o, m)},$$

where $f_s(\mathcal{M}_i, m)$ refers to the score earned by \mathcal{M}_i at the m -th match ($1 \leq m \leq K$, K is the number of performed matches):

- For zero-sum games, e.g., Tic-Tac-Toe,

$$f_s(\mathcal{M}_i, m) = \begin{cases} 1, & \text{if } \mathcal{M}_i \text{ wins at the } m\text{-th match} \\ 0, & \text{if } \mathcal{M}_i \text{ loses at the } m\text{-th match} \\ 0.5, & \text{if } \mathcal{M}_i \text{ and } \mathcal{M}_o \text{ achieve a draw} \end{cases}$$

- For non-zero-sum games, e.g., Blind Auction, $f_s(\mathcal{M}_i, m)$ is the rewards earned by \mathcal{M}_i at the m -th match.

$NRA(\mathcal{M}_i, \mathcal{M}_o, f_s)$ is naturally normalized to $[-1, 1]$, providing an interpretable meaning regarding the performance of \mathcal{M}_i : $NRA(\mathcal{M}_i, \mathcal{M}_o, f_s) > 0$ means \mathcal{M}_i is better than \mathcal{M}_o ; $NRA(\mathcal{M}_i, \mathcal{M}_o, f_s) < 0$ means \mathcal{M}_i is worse than \mathcal{M}_o ; $NRA(\mathcal{M}_i, \mathcal{M}_o, f_s) = 0$ means \mathcal{M}_i is as competitive as \mathcal{M}_o .

Evaluation Metric: Elo Rating. Following the conventional rating mechanism in the real world, e.g., Chess, we employ the popular **Elo Rating** (Elo, 1960) for calculating the relative skill levels of players in zero-sum games. Please refer to Appendix A7 for more details of Elo rating.

3 GTBENCH: Game-Theoretic Evaluation of LLMs

GTBENCH is a language-driven environment, making participating agents compete against each other in a game-theoretic manner. It is designed to be flexible and extensible, providing unified interfaces to participants and games, and supporting various multi-turn-based games which can be extended in the future. The overall framework is presented in Figure 1. There are three main components: *Environment*, *Prompt Adapter*, and *Participant*. Please refer to Appendix A1 for a detailed introduction of each component.

3.1 Taxonomy of Game-Theoretic Tasks

The chosen tasks and their detailed configurations are presented in Table 1. To comply with the common taxonomy (Lanctot et al., 2019) of game-theoretic tasks and provide diverse gaming scenarios, GTBENCH supports 10 different gaming environments, including Tic-Tac-Toe, Connect-4, Kuhn Poker, Breakthrough, Liar’s Dice, Blind Auction, Negotiation, Nim, Pig, Iterated Prisoner’s Dilemma, covering 6 mainstream game-theoretic configurations, including *complete-* and *incomplete-information* gaming, *dynamic* and *static* gaming, and *probabilistic* and *deterministic* gaming. The preferred abilities of each game could be characterized as the combination of *board strategy*, *bids*, *collaboration*, *bluff*, and *math*. Please refer to Appendix A2.1 for the rules of each game and Appendix A2.2 for an explanation of game-theoretic taxonomy.

3.2 Participants and Protocols

Conventional Agents output actions through a conventional optimization or searching process. To provide fair comparisons, we employ the powerful Monte Carol Tree Search (MCTS) (Chaslot et al., 2008) as the conventional agent for most of the games, with the number of simulations as 1000. Since Iterated Prisoner’s Dilemma is dynamic gaming with very limited action space, i.e., <TESTIFY> or <SILENT>, we utilize the more popular Tit-for-Tat (Axelrod, 1981) strategy, which simply repeating the opponent’s last action, as the conventional agent. We also include Random Agent that randomly selects action at each turn, serving as a baseline and sanity check. Please refer to Appendix A3.1 for more details about MCTS Agent and Tit-for-Tat Agent.

LLM-Driven Reasoning Agent consists of backbone LLMs and reasoning paradigms. For reasoning schemes, we consider the following reasoning paradigms as they are widely known to be effective for general reasoning tasks: ❶ *Prompt*: Directly Prompt LLMs to generate responses, without additional reasoning steps; ❷ *Chain-of-Thought (CoT)* (Wei et al., 2022): CoT Agent prompts LLMs by thinking step by step; ❸ *Self-Consistent CoT* (Wang et al., 2022b): SC-CoT Agent prompts LLMs by generating multiple step-by-step thinking trajectories and performing majority voting to get the final response. The number of trajectories is set to 5 in this paper; ❹ *Tree-of-Thought (ToT)* (Yao et al., 2024a): ToT Agent prompts LLMs to generate responses by incorporating exploration and deliberate decision-making, e.g., self-evaluation. The number of sequences for both answer generation and answer evaluations is set to 3.

Prompt Templates. Prompts are designed to be modular, consisting of four individual components: *System Prompt*, *Head Prompt*, *Observation Prompt*, and *Reasoning Prompt*. Reasoning prompts, e.g., CoT/ToT, are designed to only focus on instructing LLM how to think, regardless of the game environment. Thus, they could be automatically adapted when adding a new game. Please refer to Appendix A5 for the detailed prompts and observations for each game and agent.

Sanity Check. We provide the task completion rates of all the LLMs and reasoning agents in Appendix A5.6. We show that all the LLM agents achieve $\geq 90\%$ completion rate, indicating that the prompts are properly configured and LLMs are capable of following instructions to finish the game.

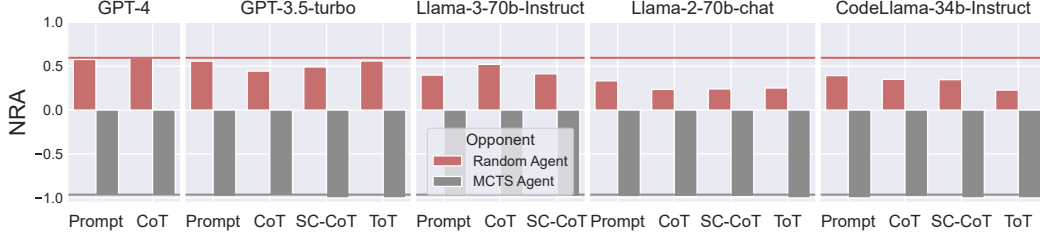


Figure 2: The NRA of state-of-the-art LLM-driven reasoning agents when against MCTS Agents and Random Agents, over complete and deterministic scenarios. Red and gray lines mean the maximum NRA achieved by LLM agents.

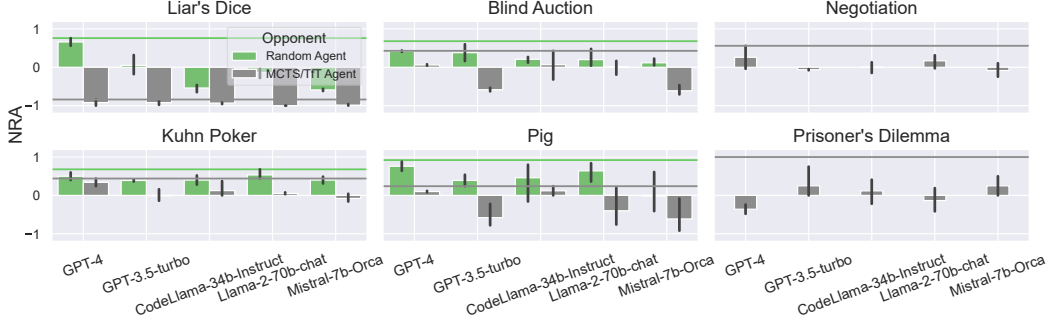


Figure 3: The game-wise NRA of LLMs when against MCTS/TfT Agents and Random Agents, over incomplete and probabilistic scenarios. Error bars are obtained over different reasoning methods. Green and gray lines mean the maximum NRA achieved by LLM agents.

4 Are LLMs Capable of Strategic Reasoning?

In this section, we evaluate the strategic reasoning capabilities of LLMs by conducting experiments among conventional solvers and LLM-driven agents.

Experimental Settings. We consider well-recognized LLMs such as commercial LLMs: GPT-3.5-turbo-1106 and GPT-4-0613 (Achiam et al., 2023), and open-source LLMs: Llama-3-70b-Instruct (Meta, 2024), Deepseek-LLM-67b-chat (Bi et al., 2024), Llama-2-70b-chat (Touvron et al., 2023), CodeLlama (Roziere et al., 2023), and Mistral-7b-Orca (Jiang et al., 2023a; Mukherjee et al., 2023). For all the LLMs, the temperature is set to 0.2 and the max number of generated tokens is 1024. For each competition, we run 50 valid matches. The final performance is measured by the averaged NRA over the 50 valid matches. To mitigate the first-player advantage, we have each participant take the first turn in 25 matches.

4.1 Complete and Deterministic Gaming

There are four complete and deterministic tasks supported in GTBENCH: Tic-Tac-Toe, Connect-4, Breakthrough, and Nim. We compare LLM-driven agents with Random Agent and MCTS Agent. Results are summarized in Figure 2. In general, we show that all LLMs achieve substantial relative advantages when competing against the Random Agent. Among all the agents, GPT-4 w/ CoT reasoning achieves the highest NRA. For open-source LLMs, Llama-3-70b-Instruct outperforms other open-source LLMs, achieving comparable capabilities as GPT-4.

However, when competing against the MCTS Agent, all the LLM agents equipped with various reasoning methods achieve NRA as -1 , meaning that LLM agents can barely win even a single match. This is because for board games with moderate action/state space such as the four involved complete and deterministic games in GTBENCH, MCTS agents with a sufficient number of simulations can achieve near-optimal strategies. Consequently, LLMs are not competitive in complete and deterministic games.

4.2 Probabilistic and Dynamic Gaming

There are five probabilistic game-theoretic gaming tasks: Kuhn Poker, Liar’s Dice, Blind Auction, Negotiation, Pig, and one dynamic task: Iterated Prisoner’s Dilemma. We group these games together as they all involve stochasticity in the gameplay, which is essentially different from complete and deterministic games. The Random Agent as the opponent is omitted for both Negotiation and Iterated Prisoner’s Dilemma because the Random Agent rarely chooses to collaborate, resulting in meaningless evaluation. Results are summarized in Figure 3. When competing against the MCTS Agent, it is shown that Liar’s Dice shares a similar trend as the complete and deterministic scenarios (Figure 2), where LLM-driven agents achieve near -1 NRA. This is because the 2-player Liar’s Dice has very limited stochasticity, making the gameplay tend to be complete information. For other tasks, we found that LLMs do not always fail. We observe that the NRA of LLM agents is close to 0 over all the tasks, indicating that they are equally competitive as conventional solvers or even better (e.g., Kuhn Poker where GPT-4 outperforms MCTS Agent).

4.3 LLM-vs.-LLM Competition

We investigate whether popular LLMs remain competitive in game-theoretic scenarios. Specifically, we take GPT-3.5-turbo with Prompt Agent as the *common opponent* and make other LLM-driven agents compete against it. Please refer to Figure A6 for the full leaderboard evaluated by NRA. The Elo rating results are placed in Table 6. In general, GPT-4 is the most powerful LLM in strategic reasoning among all the examined LLMs. Moreover, Llama-3-70b-Instruct achieves comparable performances as GPT-4 and outperforms GPT-3.5-turbo. Here we break the results into 3 takeaways:

Code-Pretraining Benefits Game-Theoretic Tasks. In Table 2, we show code-pretrained LLMs, e.g., CodeLlama-34b-Instruct and Deepseek-Coder-6.7b-Instruct, significantly outperform larger chat LLMs, e.g., Llama-2-70b-chat and Deepseek-LLM-67b-chat. These code-pretrained LLMs have less than half of the parameters, suggesting that code-pretraining benefits game-theoretic tasks. This verifies recent discoveries where code-pretraining benefits logical reasoning (Madaan et al., 2022; Liang et al., 2022; Ma et al., 2023).

Advanced Reasoning Methods Do Not Always Help. We observe that advanced reasoning methods may lead to worse results in game-theoretic scenarios. To make it more clear, we present the averaged NRA obtained by reasoning methods across different LLMs when against Random Agent in Figure 4. In general, only Mistral-7b-Orca has a substantial improvement when equipped with CoT reasoning while advanced reasoning leads to worse results for other LLMs.

In Table 3, we present the results when against GPT-3.5-turbo w/ Prompt Agent. We show that advanced reasoning benefits powerful LLMs, e.g., GPT-3.5-turbo, while it results in worse results for other LLMs. It suggests that advanced reasoning is a

Table 2: Code-pretraining benefits strategic reasoning. Gray rows are code-pretrained LLMs.

Model	avg. NRA in Det. Games	avg. NRA in Prob.	avg. NRA
GPT-4	0.09	0.15	0.13
Llama-3-70b-Instruct	-0.07	0.11	0.04
Llama-2-70b-chat	-0.25	-0.17	-0.20
CodeLlama-34b-Instruct	-0.05	0.02	-0.01
Deepseek-LLM-7b-chat	-0.09	-0.08	-0.08
Deepseek-LLM-67b-chat	0.10	-0.17	-0.05
Deepseek-Coder-6.7b-instruct	-0.14	0.07	-0.03

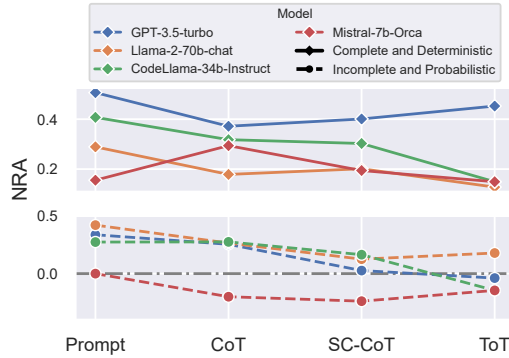


Figure 4: The NRA of LLM agents when competing against Random Agent. Advanced reasoning does not always result in better results.

Table 3: The NRA of LLM agents w/ CoT reasoning. Cyan cells mean CoT results in better performance. Magenta cells mean CoT results in worse performance.

Opponent	Model	Reasoning	avg. NRA ↑
GPT-3.5-turbo w/ Prompt Agent	GPT-3.5-turbo	Prompt	0.00
		CoT	0.02
	Llama-3-70b-Instruct	Prompt	0.04
		CoT	0.07
GPT-4 w/ Prompt Agent	GPT-4	Prompt	0.13
		CoT	0.13
	CodeLlama-34b-Instruct	Prompt	-0.01
		CoT	-0.09
GPT-4 w/ Prompt Agent	CodeLlama-34b-Instruct	Prompt	-0.01
		CoT	-0.04
	Llama-2-70b-chat	Prompt	-0.10
		CoT	-0.23

double-edged sword: ❶ powerful LLMs are capable of leveraging advanced reasoning to achieve better results; ❷ advanced reasoning may also impose reasoning errors and risks during the inference of ordinary LLMs. In Appendix A8, we further examine five different CoT strategies over the GPT-3.5-turbo model to mitigate the effect brought by prompt sensitivity, along with some failure cases presented. These CoT prompts resulting in different performances are all worse than the naive Prompt Agent.

Most Open-source LLMs are Less Competitive than Commercial LLMs in Complex Games.

We observe that most of open-source LLMs such as Llama-2-70b-chat and CodeLlama-34b-Instruct are not good at games with complex rules and board states. In Table 4, we present the average NRA when including and excluding Breakthrough³. It is shown that both Llama-2-70b-chat and CodeLlama-34b-Instruct fail in Breakthrough, resulting in worse NRA scores than GPT-4. However, we found that the recently released Llama-3-70b-Instruct (Meta, 2024) has a significant performance in Breakthrough. This indicates that open-source LLMs achieve comparable capabilities when dealing with complex tasks and environments as commercial LLMs.

Table 4: The average NRA of LLM-driven agents when Breakthrough is included and excluded.

Taxonomy	GPT-4	Llama-3-70b-Instruct	CodeLlama-34b-Instruct	Llama-2-70b-chat
w Breakthrough	0.13	0.04	-0.01	-0.20
w/o Breakthrough	0.11 (-0.02)	-0.01 (-0.05)	0.08 (+0.09)	-0.18 (+0.02)

4.4 Error Profiles

We introduce the most prevalent mistake patterns observed across different games, comprising *Misinterpretation*, *Factual Inaccuracies*, *Overconfidence*, *Calculation Mistakes*, and *Endgame*:

Misinterpretation denotes the misinterpretation of the game’s current state by LLMs, including errors like misattributing piece ownership and failing to recognize vacant spots on the board. **Factual Errors** refer to situations where the

Table 5: Quantitative results of error patterns.

Model	Percentage of Error Patterns (%)				
	Endgame Mis-detection	Mis-interpretation	Over-confidence	Calculation Error	Factual Error
GPT-4	33.33	9.80	15.69	9.80	45.10

player has a reasonable plan but their actions do not align with their plan. For instance, in Breakthrough, GPT-4 w/ CoT agent plans to fend off frontal attacks by the opponent, which is reasonable. However, it takes rear pieces to achieve that, which is impossible. **Overconfidence** describes a scenario where a player overlooks potential risks in pursuit of greater rewards. **Calculation Errors** refer to errors that occur in arithmetic, such as calculating XOR in Nim. **Endgame Mis-detection** means a failure to recognize immediate win/lose situations, e.g., a player fails to recognize a potential winning move. Demonstrations of each mistake pattern are presented in Appendix A9.

In Table 5, we present the quantitative results regarding these error patterns. It is obtained from GPT-4 w/ CoT agent when playing against conventional solvers, e.g., MCTS/TfT agent, as the opponent. We manually examined a total of 157 turns (50 matches, with 5 turns per match). We observe that LLM agents are capable of generating reasonable planning/strategies. However, they have difficulties in selecting the correct actions to align with their thoughts. Also, LLMs miss endgame situations, leading to a failure to recognize winning and losing moves.

Table 6: The Elo rating results of LLM-vs.-LLM experiments.

Model	Tic-Tac-Toe	Breakthrough	Blind Auction	Kuhn Poker	Liar’s Dice	avg. Elo
GPT-4	1554.34	1667.11	1581.94	1479.87	1676.70	1591.99
Llama-3-70b-Instruct	1371.68	1669.42	1524.11	1625.46	1694.64	1577.06
GPT-3.5-turbo	1579.80	1576.37	1514.27	1441.80	1459.26	1514.30
CodeLlama-34b-Instruct	1589.94	1398.10	1533.48	1414.57	1374.40	1462.10
Llama-2-70b-chat	1479.08	1320.42	1484.32	1521.82	1485.00	1458.13
Mistral-7B-Instruct	1440.15	1338.57	1361.89	1516.48	1310.00	1393.42

³Breakthrough has larger action/state space than other complete-information games.

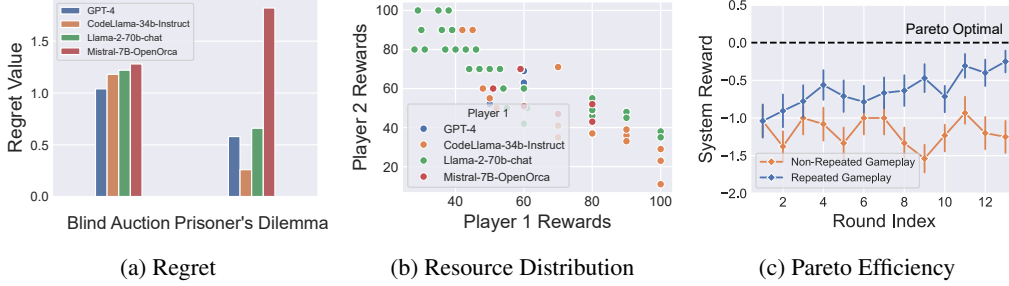


Figure 5: Game-theoretic properties. The results are obtained when competing against GPT-3.5-turbo w/ Prompt Agent as the opponent. In (b), each dot (x, y) represents an agreement in a resource distribution with Player 1 obtaining reward x and Player 2 obtaining reward y . In (c), the system reward is calculated by the sum of the payoffs of all players.

5 The Game-Theoretic Properties of LLMs

Nash Equilibrium with Regret. In game theory, being close to a Nash Equilibrium (Nash Jr, 1950) indicates that the strategies chosen by the players are near to optimal. It has been popular to approximate Nash Equilibrium with Regret⁴ (Johanson et al., 2012; Nisan and Noti, 2017; Zinkevich et al., 2007). In Figure 5a, we present the regret values of LLMs on Blind Auction and Iterated Prisoner’s Dilemma. Please refer to Appendix A10 for how regret values are calculated for these two tasks. For Blind Auction, GPT-4 shows lower Regret, indicating achieving closer to optimal solutions than other LLMs. However, in Iterated Prisoner’s Dilemma, CodeLlama-34b-Instruct exhibits lower regret compared to GPT-4. Through human examination, we found that this is because GPT-4 tends to `<Silent>` more frequently, whereas CodeLlama has a significantly higher probability of `<Testify>`. This discrepancy may be due to the human preference alignment in GPT-4, such as a higher emphasis on morality (Pan et al., 2023) or maximizing system reward⁵, which makes GPT-4 less likely to `<Testify>`.

Pareto Efficiency. We study Pareto Efficiency in two games: Negotiation and (Iterated) Prisoner’s Dilemma. In Figure 5b, we count all agreements reached by participants and record the values attributed to each based on the agreed division. Most agreements result in substantial values for both participants, though some LLMs, like Llama-2-70b-chat and CodeLlama-34b-Instruct, may accept unfair resource divisions. In contrast, GPT-4 and Mistral struggle to reach agreements and tend to negotiate for Pareto improvements. A repeated game is a standard game that is played multiple times by the same players, with each player is able to observe the history of past plays (Aumann et al., 1995; Akata et al., 2023). In Figure 5c, we investigate the Pareto Improvement in Iterated Prisoner’s Dilemma and ordinary Prisoner’s Dilemma, i.e., each round is played individually. The Pareto Improvement is observed in the repeated-game scenario during the rounds, indicating that LLMs are capable of leveraging history to adjust their strategies.

6 Conclusion

This work investigated LLMs’ strategic and logical reasoning abilities under competitive scenarios. To achieve this, we created a broad evaluation scope by considering various classic and LLM-based gaming agents and 10 representative games. We conducted the benchmark study of game-theoretic evaluations for LLMs, shedding light on their reasoning performance. Our extensive evaluations revealed insightful LLMs’ gaming behavior, such as their intrinsic failure in complete and deterministic games, impressive reasoning in incomplete and probabilistic games, and benefiting from code-generation pertaining and appropriate prompt designs.

Limitations This research prompts LLMs to generate actions regarding various game scenarios, relying on pre-defined prompt templates. Thus, the results may suffer from certain variances introduced by prompt sensitivities. Although the introduced games are popular, their actions/state

⁴Regret (Zinkevich et al., 2007) measures how much a player would have improved their outcome by choosing a different strategy, given what they know now after the game has played out.

⁵`<Silent>` maximizes the system reward in Iterated Prisoner’s Dilemma

space is limited, which may not be well-distinguished for LLMs in the same skill levels. The generated actions may be illegal due to the incapacabilities of the following instructions.

Impact Statements This paper examines the game-theoretic task proficiency of AI models. We acknowledge concerns about models becoming autonomous entities with their own objectives, especially in deception or negotiation scenarios. It’s important to note that our research measures the current capabilities of models, rather than enhancing their abilities. We do not train AI models to be competent in game theory tasks or to bluff or defect. Instead, we assess existing competencies, contributing to a deeper understanding that can inform innovative measures against potential risks. We believe our work paves the way for responsible and effective AI safety.

Acknowledgement

This work was performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under Contract No. DE- AC52-07NA27344 and was supported by the LLNL LDRD Program under Project No. 23-ERD-030 and 24-ERD-058. This work was partially supported by the NSF award FMITF-2319242. It was also partially supported by NSF-AI Engage Institute DRL-2112635 and DARPA MCS Grant N66001-19-2-4031. The views, opinions, and/or findings contained in this article are those of the authors and not of the funding agency.

References

- Sahar Abdelnabi, Amr Gomaa, Sarath Sivaprasad, Lea Schönherr, and Mario Fritz. Llm-deliberation: Evaluating llms with interactive multi-agent negotiation games, 2023.
- Marwa Abdulhai, Isadora White, Charlie Snell, Charles Sun, Joey Hong, Yuexiang Zhai, Kelvin Xu, and Sergey Levine. Lmrl gym: Benchmarks for multi-turn reinforcement learning with language models, 2023.
- 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 arXiv:2303.08774*, 2023.
- Ehsan Aghaei, Xi Niu, Waseem Shadid, and Ehab Al-Shaer. Securebert: A domain-specific language model for cybersecurity. In *International Conference on Security and Privacy in Communication Systems*, pages 39–56. Springer, 2022.
- Elif Akata, Lion Schulz, Julian Coda-Forno, Seong Joon Oh, Matthias Bethge, and Eric Schulz. Playing repeated games with large language models, 2023.
- Kimia Ameri, Michael Hempel, Hamid Sharif, Juan Lopez Jr, and Kalyan Perumalla. Cybert: Cybersecurity claim classification by fine-tuning the bert language model. *Journal of Cybersecurity and Privacy*, 1(4):615–637, 2021.
- Robert J Aumann, Michael Maschler, and Richard E Stearns. *Repeated games with incomplete information*. MIT press, 1995.
- Robert Axelrod. The emergence of cooperation among egoists. *American political science review*, 75 (2):306–318, 1981.
- Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, Athul Paul Jacob, Mojtaba Komeili, Karthik Konath, Minae Kwon, Adam Lerer, Mike Lewis, Alexander H. Miller, Sandra Mitts, Adithya Renduchintala, Stephen Roller, Dirk Rowe, Weiyan Shi, Joe Spisak, Alexander Wei, David J. Wu, Hugh Zhang, and Markus Zijlstra. Human-level play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378:1067 – 1074, 2022.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687*, 2023.

- Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*, 2024.
- Kranti Chalamalasetti, Jana Götze, Sherzod Hakimov, Brielen Madureira, Philipp Sadler, and David Schlangen. clembench: Using game play to evaluate chat-optimized language models as conversational agents. *arXiv preprint arXiv:2305.13455*, 2023.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*, 2023.
- Guillaume Chaslot, Sander Bakkes, Istvan Szita, and Pieter Spronck. Monte-carlo tree search: A new framework for game ai. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, pages 216–217, 2008.
- Junzhe Chen, Xuming Hu, Shuodi Liu, Shiyu Huang, Wei-Wei Tu, Zhaofeng He, and Lijie Wen. Llmarena: Assessing capabilities of large language models in dynamic multi-agent environments. *arXiv preprint arXiv:2402.16499*, 2024.
- Justin Chih-Yao Chen, Swarnadeep Saha, and Mohit Bansal. Reconcile: Round-table conference improves reasoning via consensus among diverse llms. *arXiv preprint arXiv:2309.13007*, 2023.
- Marc-Alexandre Côté, Akos Kádár, Xingdi Yuan, Ben Kybartas, Tavian Barnes, Emery Fine, James Moore, Matthew Hausknecht, Layla El Asri, Mahmoud Adada, et al. Textworld: A learning environment for text-based games. In *Computer Games: 7th Workshop, CGW 2018, Held in Conjunction with the 27th International Conference on Artificial Intelligence, IJCAI 2018, Stockholm, Sweden, July 13, 2018, Revised Selected Papers 7*, pages 41–75. Springer, 2019.
- Tim Ruben Davidson, Veniamin Veselovsky, Michal Kosinski, and Robert West. Evaluating language models through negotiations. In *The Twelfth International Conference on Learning Representations*, 2023.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*, 2023.
- Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5050–5063, 2024a.
- 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)*, pages 2232–2246, 2024b.
- Arpad Elo. Elo rating system. https://en.wikipedia.org/wiki/Elo_rating_system, 1960.
- Meta Fundamental AI Research Diplomacy Team (FAIR)[†], Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, Athul Paul Jacob, Mojtaba Komeili, Karthik Konath, Minae Kwon, Adam Lerer, Mike Lewis, Alexander H. Miller, Sasha Mitts, Adithya Renduchintala, Stephen Roller, Dirk Rowe, Weiyan Shi, Joe Spisak, Alexander Wei, David Wu, Hugh Zhang, and Markus Zijlstra. Human-level play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378(6624): 1067–1074, 2022.
- Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35: 18343–18362, 2022.

- Kanishk Gandhi, Dorsa Sadigh, and Noah D. Goodman. Strategic reasoning with language models, 2023.
- Ran Gong, Qiuyuan Huang, Xiaojian Ma, Hoi Vo, Zane Durante, Yusuke Noda, Zilong Zheng, Song-Chun Zhu, Demetri Terzopoulos, Li Fei-Fei, and Jianfeng Gao. Mindagent: Emergent gaming interaction, 2023.
- Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Xingdi Yuan. Interactive fiction games: A colossal adventure. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7903–7910, 2020.
- 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 arXiv:2310.06825*, 2023a.
- Haitao Jiang, Lin Ge, Yuhe Gao, Jianian Wang, and Rui Song. Large language model for causal decision making. *arXiv preprint arXiv:2312.17122*, 2023b.
- Michael Johanson, Nolan Bard, Marc Lanctot, Richard G Gibson, and Michael Bowling. Efficient nash equilibrium approximation through monte carlo counterfactual regret minimization. In *Aamas*, pages 837–846, 2012.
- Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Satyaki Upadhyay, Julien Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, et al. Openspiel: A framework for reinforcement learning in games. *arXiv preprint arXiv:1908.09453*, 2019.
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for "mind" exploration of large language model society. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*, 2022.
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and Shuming Shi. Encouraging divergent thinking in large language models through multi-agent debate. *arXiv preprint arXiv:2305.19118*, 2023.
- Jonathan Light, Min Cai, Sheng Shen, and Ziniu Hu. Avalonbench: Evaluating llms playing the game of avalon. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. Agentbench: Evaluating llms as agents, 2023.
- Yingwei Ma, Yue Liu, Yue Yu, Yuanliang Zhang, Yu Jiang, Changjian Wang, and Shanshan Li. At which training stage does code data help llms reasoning? *arXiv preprint arXiv:2309.16298*, 2023.
- Aman Madaan, Shuyan Zhou, Uri Alon, Yiming Yang, and Graham Neubig. Language models of code are few-shot commonsense learners. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1384–1403, 2022.
- Meta. Introducing meta llama 3: The most capable openly available llm to date. <https://ai.meta.com/blog/meta-llama-3/>, 2024. Accessed: 2024-05-18.
- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4, 2023.
- Gabriel Mukobi, Hannah Erlebach, Niklas Lauffer, Lewis Hammond, Alan Chan, and Jesse Clifton. Welfare diplomacy: Benchmarking language model cooperation. In *Socially Responsible Language Modelling Research*, 2023.
- John F Nash Jr. Equilibrium points in n-person games. *Proceedings of the national academy of sciences*, 36(1):48–49, 1950.

- Noam Nisan and Gali Noti. An experimental evaluation of regret-based econometrics. In *Proceedings of the 26th International Conference on World Wide Web*, pages 73–81, 2017.
- Aidan O’Gara. Hoodwinked: Deception and cooperation in a text-based game for language models, 2023.
- Alexander Pan, Jun Shern Chan, Andy Zou, Nathaniel Li, Steven Basart, Thomas Woodside, Hanlin Zhang, Scott Emmons, and Dan Hendrycks. Do the rewards justify the means? measuring trade-offs between rewards and ethical behavior in the machiavelli benchmark. In *International Conference on Machine Learning*, pages 26837–26867. PMLR, 2023.
- Rajkumar Ramamurthy, Prithviraj Ammanabrolu, Kianté Brantley, Jack Hessel, Rafet Sifa, Christian Bauckhage, Hannaneh Hajishirzi, and Yejin Choi. Is reinforcement learning (not) for natural language processing: Benchmarks, baselines, and building blocks for natural language policy optimization. In *The Eleventh International Conference on Learning Representations*, 2023.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Cote, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. In *International Conference on Learning Representations*, 2020.
- Simon Stepputtis, Joseph Campbell, Yaqi Xie, Zhengyang Qi, Wenxin Sharon Zhang, Ruiyi Wang, Sanketh Rangreji, Charles Michael Lewis, and Katia P Sycara. Long-horizon dialogue understanding for role identification in the game of avalon with large language models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023a.
- Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. Scienceworld: Is your agent smarter than a 5th grader? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11279–11298, 2022a.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2022b.
- Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. Mint: Evaluating llms in multi-turn interaction with tools and language feedback. *arXiv preprint arXiv:2309.10691*, 2023b.
- 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.
- Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabrovolski, Mark Dredze, Sebastian Gehrmann, Prabhajan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564*, 2023.
- Kai Xiong, Xiao Ding, Yixin Cao, Ting Liu, and Bing Qin. Examining the inter-consistency of large language models: An in-depth analysis via debate. Association for Computational Linguistics, 2023.
- 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 arXiv:2309.04658*, 2023a.

- Zelai Xu, Chao Yu, Fei Fang, Yu Wang, and Yi Wu. Language agents with reinforcement learning for strategic play in the werewolf game, 2023b.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757, 2022.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024a.
- Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, page 100211, 2024b.
- Yadong Zhang, Shaoguang Mao, Tao Ge, Xun Wang, Adrian de Wynter, Yan Xia, Wenshan Wu, Ting Song, Man Lan, and Furu Wei. Llm as a mastermind: A survey of strategic reasoning with large language models. *arXiv preprint arXiv:2404.01230*, 2024.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.
- Andrew Zhu, Karmanya Aggarwal, Alexander Feng, Lara Martin, and Chris Callison-Burch. Fireball: A dataset of dungeons and dragons actual-play with structured game state information. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2023.
- Martin Zinkevich, Michael Johanson, Michael Bowling, and Carmelo Piccione. Regret minimization in games with incomplete information. *Advances in neural information processing systems*, 20, 2007.

Appendix

A1 Overall Architecture

There are three main components in GTBENCH:

- **Environment.** The environment (Figure 1 (c)) is responsible for overseeing the crucial processes related to gameplay. Specifically, it is tasked with building up observations, managing gameplay, and applying the actions obtained from participants. In this paper, all of the gaming environments are built on top of OpenSpiel [Lanctot et al. \(2019\)](#).
- **Prompt Adapter.** The prompt adapter (Figure 1 (b)) plays a vital role in facilitating effective communication between the environment and the virtual participants. It serves as an intermediary between the two entities by receiving observations from the environment, which it then translates into unified observation prompts. The prompts are then parsed and sent to the participating agents to formulate their responses. The adapter is also responsible for obtaining actions from the participants, which it transforms into legal actions before parsing them to the environment for game execution.
- **Participant.** The participants (Figure 1 (a)) involved in the gaming process generate responses according to the observation prompts received from the Prompt Adapter. These responses consist of actions that participants intend to take in this turn.

A2 Gameplay Configurations

A2.1 Games Introduction

Tic-Tac-Toe⁶ is a paper-and-pencil game for two players who take turns marking the spaces in a three-by-three grid with X or O. The player who succeeds in placing three of their marks in a horizontal, vertical, or diagonal row is the winner. It is a solved game, with a forced draw assuming optimal play from both players.

- **Observation (input):** Our observation contains “opponent moves” and “self moves”. “Opponent moves” contains all the current opponent agent’s historical actions. “Self moves” contains all the current agent’s history actions.
- **Actions:** We define our action in the following format: $\langle CxRy \rangle$, in which C and R mean columns and rows respectively, while x and y mean the index of column and row. Each player may make their own action in turn.

Prisoner’s Dilemma⁷ is a game theory thought experiment that involves two rational agents, each of whom can cooperate for mutual benefit or betray their partner (“defect”) for individual reward.

- **Observation (input):** Our observation contains “opponent moves” and “self moves”. “Opponent moves” contains all the current opponent agent’s historical actions. “Self moves” contains all the current agent’s history actions.
- **Actions:** We define our action in the following format: $\langle \text{Silent} \rangle$ or $\langle \text{Testify} \rangle$. All players must take their action simultaneously.

Breakthrough⁸ Breakthrough is an abstract strategy board game invented by Dan Troyka in 2000 and made available as a Zillions of Games file (ZRF). It won the 2001 8x8 Game Design Competition. The first player to reach the opponent’s home row — the one farthest from the player — is the winner. In our work, we scale the size of the board to 3*8 while maintaining its competitiveness.

- **Observation (input):** Our observation contains “opponent moves”, “self moves”, and “board preview”. “Opponent moves” contains all the current opponent agent’s historical actions. “Self moves” contains all the current agent’s history actions. The “board preview” feature maintains the status of each grid on the board through a list of strings, denoting whether it contains a black piece, a white piece, or is empty.

⁶<https://en.wikipedia.org/wiki/Tic-tac-toe>

⁷https://en.wikipedia.org/wiki/Prisoner%27s_dilemma

⁸[https://en.wikipedia.org/wiki/Breakthrough_\(board_game\)](https://en.wikipedia.org/wiki/Breakthrough_(board_game))

- **Actions:** We define our action in the following format: $Ax \rightarrow By$, in which A and B mean the current column index and destination column index respectively, while x and y mean the index of current row and destination row. Each player may make their own action in turn.

Connect Four⁹ is a game in which the players choose a color and then take turns dropping colored tokens into a six-row, seven-column vertically suspended grid. The pieces fall straight down, occupying the lowest available space within the column. The objective of the game is to be the first to form a horizontal, vertical, or diagonal line of four of one's own tokens.

- **Observation (input):** Our observation contains “opponent moves” and “self moves”. “Opponent moves” contains all the current opponent agent’s historical actions. “Self moves” contains all the current agent’s history actions.
- **Actions:** We define our action in the following format: $\langle Cx \rangle$ in which C means column, while x means the index of column. Each player may make their action in turn.

Blind Auction¹⁰ is a common type of auction. In this type of auction, all bidders simultaneously submit sealed bids so that no bidder knows the bid of any other participant. The highest bidder pays the price that was submitted. All players must take their action simultaneously.

- **Observation (input):** Our observation contains “valuation”. “Valuation” contains each of the values of all the items for the current player.
- **Actions:** We define our action in the following format: $\langle x \rangle$, in which x represents the amount that a certain player would like to bid for.

Kuhn Poker¹¹ is a simplified form of poker. Kuhn is a simple model zero-sum two-player imperfect-information game, amenable to a complete game-theoretic analysis. In Kuhn poker, the deck includes only three playing cards, for example, a King, Queen, and Jack. One card is dealt to each player, which may place bets similarly to a standard poker. If both players bet or both players pass, the player with the higher card wins, otherwise, the betting player wins.

- **Observation (input):** Our observation contains “card”, and “moves”. Among these, “card” denotes the current player’s hand card in this match, while “moves” represents the history of all characters’ moves together with the index of the rounds.
- **Actions:** We define our action in the following format: $\langle \text{Pass} \rangle$ or $\langle \text{Bet} \rangle$. Each player may make their own action in turn.

Liar’s Dice¹² is a class of dice games for two or more players requiring the ability to deceive and detect an opponent’s deception.

- **Observation (input):** Our observation contains: “Self dice face value” and “last move”. “Self dice face value” describes all the face values of dices the current player has, while “last move” represents the previous player’s action.
- **Actions:** We define our action in the following format: $\langle x \text{ dices}, y \text{ value} \rangle$ or $\langle \text{Liar} \rangle$. Among these, x means the quantity of dice, and y means the face values of the dice. The option “Liar” denotes the current player wants to stop and challenge the previous players. Each player may make their own action in turn.

Pig¹³ is a simple dice game. Players take turns to roll a single dice as many times as they wish, adding all roll results to a running total, but losing their gained score for the turn if they roll a 1.

- **Observation (input):** Our observation contains: “self current score”, “opponent current score”, and “turn total score”. “Self current score” and “opponent current score” represent the game culminated score of the current player and opponent player respectively. While “turn total score” denotes the sum of the score of the current turn.

⁹https://en.wikipedia.org/wiki/Connect_Four

¹⁰https://en.wikipedia.org/wiki/First-price_sealed-bid_auction

¹¹https://en.wikipedia.org/wiki/Kuhn_poker

¹²https://en.wikipedia.org/wiki/Liar%27s_dice

¹³[https://en.wikipedia.org/wiki/Pig_\(dice_game\)](https://en.wikipedia.org/wiki/Pig_(dice_game))

- **Actions:** We define our action in the following format: $\langle \text{stop} \rangle$ or $\langle \text{roll} \rangle$. Each player may make their own action in turn.

Nim¹⁴ is a mathematical game of strategy in which two players take turns removing objects from distinct heaps or piles. On each turn, a player must remove at least one object and may remove any number of objects provided they all come from the same heap or pile.

- **Observation (input):** Our observation contains: “piles”. “Piles” denotes the number of matches different piles have.
- **Actions:** We define our action in the following format: $\langle \text{pile}:x, \text{take}:y \rangle$. Among these, x represents the index of the pile that the current player takes, and y represents the number of matches the current player takes. Each player may make their own action in turn.

Negotiation¹⁵

- **Observation (input):** Our observation contains: “turn type”, “item pool”, “most recent proposal”, “most recent utterance”, and “self value vector”. “turn type” is an enum variable, it has two options: proposal and utterance. The “Proposal” is the turn that the current player could think about the desired quantities of the items, and the “Utterance” is the turn that the current player states the values to its opponent. “item pool” represents the quantities of all the items. “most recent proposal” and “most recent utterance” represent the opponent’s latest proposal and utterance. “self value vector” represents how much the value of the items to the current player.
- **Actions:** We define our action in the following format: $\langle \text{Agree} \rangle$ or $\langle x, y, z \rangle$. Among these, $\langle \text{Agree} \rangle$ represents the current player agreeing on the opponent’s utterance. x , y , and z represent the quantities of different items that the current player wants to get.

A2.2 Gaming-Theoretic Taxonomy

Complete and Incomplete Information One fundamental dimension along which games are classified is the level of information available to players. In *complete information* games, players possess perfect knowledge regarding the game’s structure, including the available strategies, payoffs, and the actions taken by other players. Examples of complete information games include canonical examples like chess and Tic-Tac-Toe, where all relevant information is transparent to all players throughout the game. Conversely, *incomplete information* games involve situations where players must make decisions without having full knowledge of the game’s parameters or the actions of other players. Classic examples of incomplete information games include strategic interactions in economics, such as auctions or negotiations, where players have limited knowledge about the valuations or preferences of other participants.

Dynamic and Static Another crucial dimension for classifying games is the timing of players’ decisions. In *static games*, players make decisions simultaneously, without the opportunity to observe or react to other players’ moves. Examples of static games include simultaneous-move games like the Iterated Prisoner’s Dilemma. In contrast, *dynamic games* involve sequential decision-making, where players observe previous moves before choosing their actions. Dynamic games encompass a wide range of strategic environments, from turn-based board games like chess to dynamic settings like Kuhn Poker, where players strategically make their actions based on the unfolding dynamics of the game.

Probabilistic and Deterministic Games can also be differentiated based on the role of uncertainty in decision-making. In *deterministic games*, the outcomes of players’ actions are fully determined by the game’s rules and the strategies chosen by players. Deterministic games include classic examples like chess or Tic-Tac-Toe, where each move leads to a predictable outcome based on the game’s rules and the players’ strategies. Conversely, *probabilistic games* involve randomness or uncertainty in determining outcomes. This uncertainty can stem from elements such as dice rolls, card draws. Examples of probabilistic games include games of chance like Kuhn Poker, Liar’s Dice, or Pig, where players must contend with the inherent uncertainty of probabilistic outcomes.

¹⁴<https://en.wikipedia.org/wiki/Nim>

¹⁵<https://arxiv.org/pdf/1706.05125.pdf>

A3 Participants

A3.1 Conventional Agent

MCTS [Chaslot et al. \(2008\)](#) is a heuristic search algorithm that has gained prominence in recent years, particularly in the domain of board games and decision-making under uncertainty. It is characterized by its ability to efficiently explore large search spaces by sampling potential future outcomes through Monte Carlo simulations. The algorithm iteratively builds a search tree by simulating random sequences of moves from the current game state and evaluating their outcomes through repeated simulations. By focusing computational resources on promising branches of the search tree, MCTS aims to guide the search towards regions of the game space that are more likely to lead to favorable outcomes. MCTS has demonstrated remarkable success in various domains, including games like Go, where traditional search algorithms struggle due to the game’s immense complexity and branching factor.

Tit-for-Tat [Axelrod \(1981\)](#) is a simple but powerful strategy in the realm of repeated games and social dilemmas. The strategy is based on the principle of reciprocity, where an agent initially cooperates and then mimics the opponent’s previous action in subsequent rounds. Specifically, Tit-for-Tat starts by cooperating in the first round and then replicates the opponent’s last move in each subsequent round. Despite its simplicity, Tit-for-Tat has been shown to be remarkably effective in promoting cooperation and achieving favorable outcomes in various scenarios, including Iterated Prisoner’s Dilemma and evolutionary simulations. Its success stems from its ability to balance cooperation and retaliation, fostering reciprocal behavior and encouraging cooperation among interacting agents.

A4 LLM-vs-LLM Results

In Figure A6, we present the confusion matrix of NRA when various LLM agents are against GPT-3.5-turbo and GPT-4.

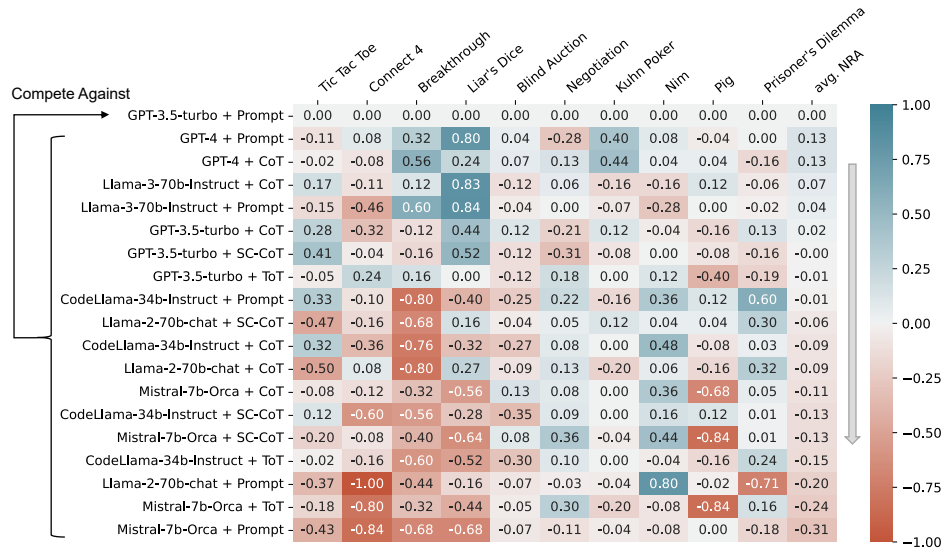


Figure A6: NRA confusion matrix of LLM vs. LLM across ten games ranked by average NRA. GPT-3.5-turbo with Prompt Agent serve as the common opponent against multiple combinations of LLMs with agents.

A5 Prompt and Protocol

A5.1 Modular Prompt Structure

When prompting LLMs to generate the next action during the course of a game, the prompt is composed of four individual components, to make sure all the participants access the same observations and information from environments:

System Prompt provides general guidance on how the LLMs should perform.

Head Prompt provides the general background and rules of the game.

Observation Prompt is formatted by a fixed game-wise template, providing sufficient observations from the environment regarding the current gaming state, to make LLMs capable of making decisions. The following provides the template used in the *Blind Auction* environment:

Your budget is <VALUATION>. Your bid must be strictly lower than or equal to <VALUATION>. Your opponent also has an expected valuation and you do not know it.
The legal actions are: <LEGAL_MOVES>.

Here <VALUATION> and <LEGAL_MOVES> are variables and are obtained from a unified <observation> object. In this way, all the participants are guaranteed to assess the same information.

Reasoning Prompt guides the LLM's generation process, e.g., "Let's think step by step" for the CoT Agent.

A5.2 System Prompt

The system prompt is the initial text or context provided by the user to prompt the model to generate a response. This prompt serves as the starting point for the model to understand the user's query or input and generate a relevant and coherent response based on the provided context. We use the same system prompt for different tasks in our work. Our system prompt is presented as follows:

System Prompt: You are a powerful gaming agent who can make proper decisions to beat the user in gaming tasks. You are a helpful assistant that strictly follows the user's instructions. You must answer your questions by choosing one of the legal moves given by the user!

A5.3 Head Prompt

The head prompt is the text that explains the rules of the games to the LLMs. In our work, we designed different prompts for each of the games.

- **Tic-Tac-Toe**

Head Prompt: Tic Tac Toe is a two-player game played on a grid. Players take turns marking a space with their respective symbols. The goal is to get 3 of one's own symbols in a row, either horizontally, vertically, or diagonally, before the opponent does. If all nine squares are filled and no player has three in a row, the game is a draw. The Tic Tac Toe game is played on a 3 by 3 grid, with the winning length as 3. Each move is represented by a string consisting of two parts: the column (C) and the row (R), in that order. For instance, C1R2 means the movement at the position of the first column and the second row of the grid. You are playing this game with the user (opponent).

- **Iterated Prisoner's Dilemma**

Head Prompt: You and your partner are in the Prisoner's Dilemma situation. Specifically, if you <Testify> against your partner and your partner remains <Silent>, you will go free while your partner will get 3 years in prison on the main charge. If you remain <Silent> but your partner <Testify> against you, you will serve 3 years in prison and your partner will be set free. If you and your partner <Testify> against each other, you and your partner will each

serve 2 years. If both you and your partner remain <Silent>, you and your partner will each serve 1 year.

- **Breakthrough**

Head Prompt: Breakthrough is a two-player game played on a rectangular board. Players take turns moving their pieces, which can move one space straight or diagonally forward if the target square is empty. A piece can also move diagonally forward to capture an opponent's piece. Capturing is optional, and a player can only capture one piece per turn. The goal is to be the first to reach the opponent's home row, the farthest row from the player. If all of a player's pieces are captured, they lose. The game does not allow draws, as pieces can only move forward or be captured. The Breakthrough board is identified by columns labeled starting from A (from left to right) and rows numbered 1 to 8 (from bottom to top). The intersection of a column and a row specifies a unique square on the board.

- **Connect Four**

Head Prompt: Connect 4 is a two-player connection board game, where the players choose a color and then take turns dropping colored discs into a vertically suspended grid. The pieces fall straight down, occupying the next available space within the column. The objective of the game is to be the first to form a horizontal, vertical, or diagonal line of four of one's own discs. You are a gaming agent who aims to beat me in Connect 4 games. Each move is represented by a string consisting of two parts: the column (C) and the row (R), in that order. For instance, C1 means the first column.

- **First-price sealed-bid auction**

Head Prompt: A first-price sealed-bid auction (FPSBA) is a common type of auction. It is also known as the blind auction. In this type of auction, all bidders simultaneously submit sealed bids so that no bidder knows the bid of any other participant. The highest bidder pays the price that was submitted.

Each action is represented by < x > where x refers to the bid.

- **Kuhn Poker**

Head Prompt: Kuhn poker is a simple model zero-sum two-player imperfect-information game, amenable to a complete game-theoretic analysis. In Kuhn poker, the deck includes only three playing cards: a King (K), a Queen (Q), and a Jack (J). One card is dealt to each player, and the third is put aside unseen. The players take turns either <Bet> to match the bet raised by the opponent or <Pass> to concede the game.

If a player bets, the other player must either call the bet by matching it or fold by conceding the game. If both players pass, the game is over, and the player with the higher-ranking card wins. The card rankings are as follows: King (K) > Queen (Q) > Jack (J).

You are playing Kuhn poker with the opponent. The actions are denoted by <Bet> and <Pass>.

- **Liar's Dice**

Head Prompt: Liar's Dice is a game of bluffing and probability, played with two players and each player has 1 dice. During each turn, a player can either bid a higher quantity of any particular face value or the same quantity of a higher face value than the previous bid. Each player tries to outbid their opponent without being caught in a lie. The move in this game is denoted in < x dices, y value>, meaning there are at least x dices with face values as y .

- **Pig**

Head Prompt: Pig is a fast-paced dice game where players risk accumulating points with each roll but risk losing them all if they roll a 1. Each player must decide when to stop rolling and bank their points, aiming to be the first to reach 100 points. You are playing Pig with the other.

- **Nim**

Head Prompt: In Nim, a strategic game with a set of four piles containing 1, 3, 5, and 7 matches respectively, players aim to avoid taking the last match. During each turn, a player may take any number of matches from a single pile, but must take at least one and cannot exceed the number remaining in that pile. The objective is to force the opponent to pick up the final match, thereby winning the game.

The action is presented in $\langle \text{pile}:x, \text{take}:y \rangle$, which means take y match(es) from the x -th pile.

- **Negotiation**

Head Prompt: You are negotiating the division of Peppers, Strawberries, and Cherries with the opponent. Different values these items hold for both you and your opponent. The process is structured into two stages per round: the proposal stage and the utterance stage.

A5.4 Observations

Our research team has developed a range of observation prompts tailored to different types of games. The list of these prompts is presented below.

- **Tic-Tac-Toe**

Observation Prompt: Your opponent has finished actions: $\langle \text{OPPONENT_MOVES} \rangle$. You have finished actions: $\langle \text{SELF_MOVES} \rangle$.

- **Iterated Prisoner's Dilemma**

Observation Prompt: You have been through this situation in the past and here are the decisions you and your partner made: (In the $idx + 1$ th round, you decided to $\langle \text{MOVE} \rangle$ and your opponent decided to $\langle \text{OPPONENT_MOVE} \rangle$) * n round

- **Breakthrough**

Observation Prompt: The board now looks like : $\langle \text{BOARD_PREVIEW} \rangle$. Among which, the letter 'b' represents a black piece, while the letter 'w' represents a white piece. And the character "." represents vacant space. The numbers in the board are the indexes of the rows. Your opponent has finished actions: $\langle \text{OPPONENT_MOVES} \rangle$. You have finished actions: $\langle \text{SELF_MOVES} \rangle$.

- **Connect Four**

Observation Prompt: Your opponent has finished actions: $\langle \text{OPPONENT_MOVES} \rangle$. You have finished actions: $\langle \text{SELF_MOVES} \rangle$.

- **First-price sealed-bid auction**

Observation Prompt: Now, you are in an auction with an opponent. You want to win the object and at the same time, your budget is $\langle \text{VALUATION} \rangle$. Your bid must be strictly lower than or equal to $\langle \text{VALUATION} \rangle$. You shall bid wisely against your opponent. Your opponent also has an expected valuation and you do not know it.

- **Kuhn Poker**

Observation Prompt: In this match, your card is $\langle \text{CARD} \rangle$. Here are the past moves in this match: $\langle \text{SELF_MOVES} \rangle$, $\langle \text{OPPONENT_MOVES} \rangle$.

- **Liar's Dice**

Observation Prompt: Currently, the face value of your dice is $\langle \text{FACE_VALUE} \rangle$. Last time, the opponent called $\langle \text{OPPONENT_LAST_ACTION} \rangle$. You are playing the Liar's Dice with another opponent. Therefore, there are only two dice in total.

- **Pig**

Observation Prompt: Right now, your current score is <AGENT_CURRENT_SCORE> and your opponent’s current score is <OPPONENT_CURRENT_SCORE>. In this turn, you have earned <TURN_TOTAL_SCORE> score.

- **Nim**

Observation Prompt: Currently, the 1st pile has <PILES[0]> match(es), the 2nd pile has <PILES[1]> match(es), the 3rd pile has <PILES[2]> match(es), 4th pile has <PILES[3]> match(es).

- **Negotiation** We proposed two different prompts for the “proposal” turn and “utterance” turn respectively.

For the “proposal” turn, we have:

Observation Prompt: Now, the opponent propose to take <OPPONENT_PROPOSAL_TAKE[0]> peppers, <OPPONENT_PROPOSAL_TAKE[1]> strawberries, and <OPPONENT_PROPOSAL_TAKE[2]> cherries from the item pool. Last time, the utterance of the opponent was to take <OPPONENT_UTTERANCE_TAKE[0]> peppers, <OPPONENT_UTTERANCE_TAKE[1]> strawberries, and <OPPONENT_UTTERANCE_TAKE[2]> cherries from the item pool.

Now, it is your decision. If you find the proposal raised by the opponent is acceptable, you should output Agree. Otherwise, you should output your proposal in the format <Proposal: [a, b, c]>.

For the “utterance” turn, we have:

Observation Prompt: Last time, you propose to take <AGENT_PROPOSAL_TAKE[0]> peppers, <AGENT_PROPOSAL_TAKE[1]> strawberries, and <AGENT_PROPOSAL_TAKE[2]> cherries from the item pool. Last time, the utterance of the opponent was to take <OPPONENT_UTTERANCE_TAKE[0]> peppers, <OPPONENT_UTTERANCE_TAKE[1]> strawberries, and <OPPONENT_UTTERANCE_TAKE[2]> cherries from the item pool.

Now, it is your turn to provide your utterance regarding the division of items. The utterance is what you want to tell your opponent and does not mean your real intent. You should output your utterance in the format <Utterance: [a, b, c]>.

A5.5 Reasoning Prompt

- **Prompt agent:** Prompt agent does not necessitate the use of LLMs to apply any predetermined strategy prior to decision-making. Rather, it simply requests LLMs for inference and subsequently provides the outcome.

You must choose a legal action to set up advantages. Your output must be in the following format:

Action: Your action wrapped with <>, i.e., <format>

Please return your answer without explanation!

- **CoT agent:** CoT agent makes LLMs consider the given observation first, then give out the action according to its thinking.

First think about your current situation, then you must choose one action from legal actions to set up advantages.

Your output must be in the following format strictly:

Thought: Your thought.

Action: Your action wrapped by <>, i.e., <format>

Remember, you can only choose one move from the legal actions.

- **SC-CoT agent:** SC-CoT agent is an advanced version of the CoT agent. It obtains actions from multiple CoT trajectories. It employs the same prompt templates as in the CoT agent.

First think about your current situation, then you must choose one action from legal actions to set up advantages.

Your output must be in the following format strictly:

Thought: Your thought.

Action: Your action wrapped by $\langle \rangle$, i.e., $\langle \text{format} \rangle$

Remember, you can only choose one move from the legal actions.

- **ToT agent:** we follow the text generation task implementation in the official codebase of ToT¹⁶. Specifically, the ToT is factorized into 1). candidate thought generation, 2). thought voting, 3). candidate action generation, 4). action voting:

Here we provide the basic prompt template used in ToT.

Step Prompt: First think about your current situation, then choose one move from legal positions to set up advantages.

Your output should be of the following format:

Thought:

Your thought.

Move:

Your action wrapped with $\langle \rangle$, e.g., $\langle \text{format} \rangle$

After executing step prompts in a breath-first search manner, we utilize the original ToT vote prompt:

Vote Prompt: Given an instruction and several choices, decide which choice is most promising. Analyze each choice in detail, then conclude in the last line "The best choice is s", where s the integer id of the choice.

A5.6 Sanity Check

To evaluate the effectiveness of our framework, we perform a sanity check by calculating the completion rates of each game. The completion rates are calculated as $\frac{50}{N}$ where N is the number of matches that will take to achieve 50 valid matches. Here, a valid match means all the participants will always generate legal moves at each turn of the match. Results are summarized in Table A7. We show that all the LLM agents achieve $\geq 90\%$ completion rate, showing that the prompts are properly configured and LLMs are capable of following instructions to finish the game.

Table A7: Sanity check. The completion rates of LLM agents over all the games.

Backbone LLM	Reasoning	Tic Tac Toe	Connect 4	Breakthrough	Liar's Dice	Blind Auction	Negotiation	Kuhn Poker	Nim	Pig	Prisoner's Dilemma	avg
GPT-3.5-turbo	Prompt	100%	100%	98%	98%	100%	100%	100%	100%	100%	100%	100%
	CoT	100%	100%	98%	100%	100%	100%	100%	98%	100%	100%	100%
	SC-CoT	100%	100%	100%	98%	100%	100%	100%	100%	100%	100%	100%
Llama-2-70b-chat	Prompt	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
	CoT	81%	98%	64%	100%	89%	69%	100%	100%	100%	98%	90%
	SC-CoT	89%	91%	81%	100%	94%	68%	100%	100%	100%	100%	92%
CodeLlama-34b-Instruct	Prompt	98%	100%	89%	100%	100%	100%	100%	100%	100%	100%	99%
	CoT	82%	100%	58%	100%	100%	78%	100%	100%	100%	100%	92%
	SC-CoT	71%	100%	71%	100%	100%	77%	100%	100%	100%	100%	92%
Mistral-7b-Orca	Prompt	98%	100%	98%	100%	100%	100%	100%	100%	100%	100%	99%
	CoT	94%	98%	100%	100%	100%	100%	100%	100%	100%	100%	99%
	SC-CoT	93%	100%	100%	100%	100%	100%	100%	100%	100%	100%	99%

A6 How Temperature Affects LLM Performance

To study how the temperature used in generating LLMs' responses affects performances, we conduct experiments by making LLMs with 0.2 temperature (the default setting as in our paper) play against LLMs with 0.4/0.6/0.8 temperature, over CodeLlama-34b-Instruct and GPT-3.5-turbo-1106. For each experiment, we run 20 matches. The reasoning method is the PromptAgent. The results are summarized as in Table A8. We show that a larger temperature will result in worse performance for deterministic games, while it has a model-specific effect for probabilistic games.

¹⁶<https://github.com/princeton-nlp/tree-of-thought-llm/blob/master/src/tot/prompts/text.py>

Table A8: The affect of various temperatures for generation sampling.

Model	Temperature	avg. NRA in Probabilistic Games	avg. NRA in Deterministic Games
CodeLlama-34b-Instruct	0.4	-0.13	-0.01
CodeLlama-34b-Instruct	0.6	-0.16	-0.05
CodeLlama-34b-Instruct	0.8	-0.16	-0.10
GPT-35-turbo	0.4	0.04	-0.10
GPT-35-turbo	0.6	0.06	-0.12
GPT-35-turbo	0.8	0.02	-0.34

A7 Elo Rating System

The Elo rating system [Elo \(1960\)](#) is a popular method for calculating the relative skill levels of players in two-player games such as Chess. It was used by various organizations to rank players. Assume there are two players: A and B , and each player has a rating, R_A, R_B , which is a numerical value representing their skill level. The expected score for a player is the probability that the player will win against another player:

$$E_A = \frac{1}{1 + 10^{(R_B - R_A)/400}}$$

$$E_B = \frac{1}{1 + 10^{(R_A - R_B)/400}}.$$

After a match between A and B , the real values, S_A and S_B , are defined as

- If Player A wins, $S_A = 1$ and $S_B = 0$
- If Player B wins, $S_A = 0$ and $S_B = 1$
- If the game is a draw, $S_A = S_B = 0.5$

Then, the updated rating R'_A and R'_B are calculated as:

$$R'_A = R_A + K * (S_A - E_A)$$

$$R'_B = R_B + K * (S_B - E_B),$$

where K is a constant that determines how much the rating changes after a game. A higher K results in a larger change. In our paper, the initial rating is set to 1500, i.e., $R_A = R_B = 1500$, and $K = 20$.

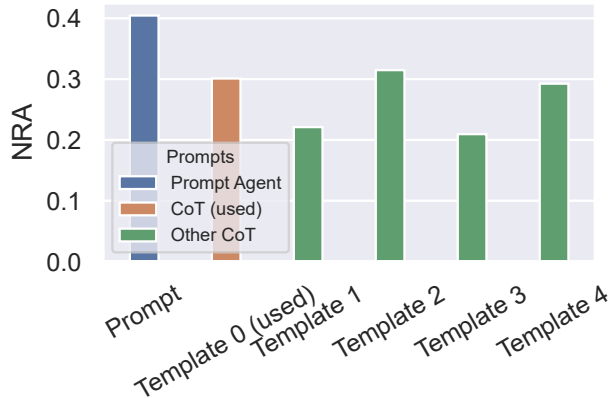


Figure A7: Investigating the sensitivity of Chain-of-Thought prompt. Prompt (used) and CoT (used) refer to the prompts utilized by the Prompt Agent and the CoT Agent in this paper. Results are obtained from the model GPT-3.5-turbo over all the game-theoretic tasks. Please refer to Table A9 for Template 0 to Template 4.

A8 Chain-of-Thought Sensitivity

We provide five different CoT strategies over the GPT-3.5-turbo model as shown in Table A9 the results presented in Figure A7.

We also include instances as shown in Table A10 where CoT agents were unable to produce legal outcomes, which serves as evidence of their limitations.

Table A9: Different Chain-of-Thought strategies.

Index	Chain-of-Thought Prompt
0	First, think about your current situation.
1	First, think about your current situation by thinking carefully, logically, and step by step.
2	First, provide your thoughts.
3	First, analyze and think about your current state and provide your strategy to beat your opponent.
4	First, analyze the state and provide your strategy to beat your opponent.

Table A10: Fail cases in different Chain-of-Thought prompts.

Index	Game	Chain-of-Thought Prompt
0	Breakthrough	Thought: I can see that my opponent has advanced their pieces quite far, and I need to make strategic moves to gain an advantage. I should prioritize moving my pieces forward while also considering potential capture opportunities. Action: <c4->b5>
1	Breakthrough	Thought: Considering the current board position and the opponent's recent moves, it seems that the best move would be to advance the white piece at b3 to c4, capturing the black piece at c4. This move would help in gaining control of the board and putting pressure on the opponent. Action: <c3->c4*>
2	Breakthrough	Thought: I see that the opponent's pieces are advancing towards my home row, and they have already captured one of my pieces. I need to start advancing my pieces and try to create opportunities for capturing the opponent's pieces to gain an advantage. Action: <b2->b3>
3	Nim	Thought: In this state, the total number of matches is 4, which is a losing position in Nim. To win, I need to force the total number of matches to be a multiple of 4 after my turn. This means I should aim to leave 3, 7, 11, or 15 matches for my opponent after my move. Action: <pile:4, take:3>
4	Nim	Thought: The current state of the game is such that the total number of matches remaining is 3, which is a losing position for the current player. To win, the strategy should be to force the opponent into a losing position by making a move that leaves 2 matches remaining. Action: <pile:4, take:6>

A9 Error Pattern Demonstrations

We present demonstrations of error patterns in Table A11.

Table A11: We conclude 5 common error patterns in GTBENCH, including misinterpretation, factual error, math calculation, and over-confidence. For each demonstration, the explanation explains why LLM generation is incorrect.

Game	Category	Demonstrations
Breakthrough	Factual Error	<p>LLM Generation: “Considering the current situation, moving the piece at b6 to either a5, b5, or c5 would allow me to advance further into the opponent’s territory. However, this could potentially expose my piece to being captured. On the other hand, moving a piece from the back row (a8, b8, or c8) to b7 could potentially block the opponent’s advance and provide a defensive line. Action: <b8->b7>”.</p> <p>Explanation: In the third round of this match, moving the back row renders it impossible to hinder the opponent’s forward movement.</p>
Nim	Calculation Error	<p>State: 3 matches in Pile 2, 5 matches in Pile 3, and 7 matches in Pile 4.</p> <p>LLM Generation: “The current state of the game is such that the 2nd pile has 3 matches, the 3rd pile has 5 matches, and the 4th pile has 7 matches. The goal is to avoid being the one to pick up the last match. To achieve this, I need to consider the binary representations of the pile sizes and aim to make the XOR of all pile sizes equal to zero. This will put me in a winning position. Currently, the XOR of all pile sizes is 4 (011 in binary). To make this zero, I need to remove 4 matches from the 3rd or 4th pile. Action: <pile:3, take:4>”.</p> <p>Explanation: The XOR of the binary representations of 3, 5, 7 is 1, not 4.</p>
Iterated Prisoner’s Dilemma	Over-confidence	<p>LLM Generation: “Given the history of the previous round, where both my partner and I chose to remain silent, it seems that we are both inclined towards cooperation. This is a positive sign, as it indicates a level of trust and mutual understanding. Considering the potential outcomes, the safest option is to remain silent. If my partner also remains silent, we both serve the least possible time. Action: <Silent>”</p> <p>Explanation: The player is overconfident that the partner will cooperate and keep silent.</p>
Connect-4	Misinterpretation	<p>LLM Generation: “I have analyzed the current situation and identified that my opponent has taken the center column, which gives them a significant advantage. To counter this, I will take the corner column to create a potential winning line. Action: <C1>.”</p> <p>Explanation: Currently, the center column is unoccupied.</p>
Tic-Tac-Toe	Endgame Misdetection	<p>State: The opponent finished actions <C2R1> and <C2R2>, the agent finished moves <C1R1> and <C3R1>.</p> <p>LLM Generation: I have one mark in the first column, and one in the third. I could potentially win if I get the chance to mark the third row of the first column in the next turn. Action: <C1R3></p> <p>Explanation: This is an immediate loss situation because the opponent has two pieces connected in a row and <C2R3> is unoccupied. The agent missed that they should take action <C2R3> to avoid lose</p>

A10 Regret Value

A10.1 Regret Value for Blind Auction

Assume that the last round of bidding was b_1 for the player and b_2 for the opponent. Assume the player's valuation is v , then the regret value is calculated by

```
if b_1 > b_2 + 1:
    regret = b_1 - (b_2 + 1)
else:
    if (b_2 + 1) < v:
        regret = v - (b_2 + 1)
    else:
        regret = 0
```

A10.2 Regret Value for Iterated Prisoner's Dilemma

The regret value Iterated Prisoner's Dilemma is simply the accumulation of the regret value of per-turn Prisoner's Dilemma:

```
if player_move == 'Testify' and opponent_move == 'Silent':
    regret = 0
elif player_move == 'Testify' and opponent_move == 'Testify':
    regret = 0
elif player_move == 'Silent' and opponent_move == 'Testify':
    regret = 1
else:
    regret = 2
```

A11 User Interfaces of GTBench Leaderboard

The user interfaces of GTBENCH leaderboard are presented in Figures [A8](#) and [A9](#).



GTBench: Uncovering the Strategic Reasoning Limitation of LLMs via Game-Theoretic Evaluations

GTBench aims to evaluate and rank LLMs' reasoning abilities in competitive environments through game-theoretic tasks, e.g., board and card games. It utilizes 10 widely recognized games supported by [OpenSpiel](#) and evaluate well-recognized LLM agents in a language-driven manner. The evaluation code and prompt templates can be found in [GTBench](#).

Please refer to [About](#) for more details of games and metrics.

The template is borrowed from [Open LLM Leaderboard](#).

GTBench

About

Select columns to show

☐ Average☐ Breakthrough☐ Connect Four☐ Blind Auction☐ Kuhn Poker☐ Liar's Dice☐ Negotiation☐ Nim☐ Pig☐ Iterated Prisoner's Dilemma☐ Tic-Tac-Toe

Model

☐ GPT-4☐ Llama-2-70b-chat☐ GPT-3.5-turbo☐ Codellama-34b-instruct☐ Mistral-7b-Orca

Agents

☐ ToT☐ Prompt☐ CoT☐ SC-CoT

Opponent Model

☐ GPT-3.5-turbo-1106☐ GPT-4

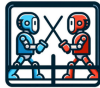
Opponent Agents

☐ prompt agent

Model	Agent	Opponent Model	Opponent Agent	Average	Tic-Tac-Toe	Connect Four	Breakthrough	Liar's Dice	Blind Auction	Ne
GPT-3.5-turbo	Prompt	GPT-3.5-turbo-1106	prompt agent	0	0	0	0	0	0	0
GPT-4	Prompt	GPT-3.5-turbo-1106	prompt agent	0.129	-0.111	0.08	0.32	0.8	0.04	-0
GPT-4	CoT	GPT-3.5-turbo-1106	prompt agent	0.126	-0.022	-0.08	0.56	0.24	0.069	0.
GPT-3.5-turbo	CoT	GPT-3.5-turbo-1106	prompt agent	0.023	0.277	-0.32	-0.12	0.44	0.115	-0
GPT-3.5-turbo	SC-CoT	GPT-3.5-turbo-1106	prompt agent	-0.002	0.409	-0.04	-0.16	0.52	-0.12	-0
GPT-3.5-turbo	ToT	GPT-3.5-turbo-1106	prompt agent	-0.005	-0.045	0.24	0.16	0	-0.12	0.
Codellama-34b-instruct	Prompt	GPT-3.5-turbo-1106	prompt agent	-0.008	0.333	-0.1	-0.8	-0.4	-0.25	0.
Llama-2-70b-chat	SC-CoT	GPT-3.5-turbo-1106	prompt agent	-0.064	-0.469	-0.16	-0.68	0.16	-0.04	0.
Codellama-34b-instruct	CoT	GPT-3.5-turbo-1106	prompt agent	-0.088	0.316	-0.36	-0.76	-0.32	-0.268	0.
Llama-2-70b-chat	CoT	GPT-3.5-turbo-1106	prompt agent	-0.089	-0.5	0.08	-0.8	0.265	-0.086	0.
Mistral-7b-Orca	CoT	GPT-3.5-turbo-1106	prompt agent	-0.113	-0.077	-0.12	-0.32	-0.56	0.133	0.
Codellama-34b-instruct	SC-CoT	GPT-3.5-turbo-1106	prompt agent	-0.128	0.122	-0.6	-0.56	-0.28	-0.348	0.

Figure A8: The user interface of GTBENCH leaderboard.

A28



GTBench: Uncovering the Strategic Reasoning Limitation of LLMs via Game-Theoretic Evaluations

GTBench aims to evaluate and rank LLMs' reasoning abilities in competitive environments through game-theoretic tasks, e.g., board and card games. It utilizes 10 widely recognized games supported by [OpenSpiel](#) and evaluate well-recognized LLM agents in a language-driven manner. The evaluation code and prompt templates can be found in [GTBench](#).

Please refer to [About](#) for more details of games and metrics.

The template is borrowed from [Open LLM Leaderboard](#).

GTBench

About

Select columns to show

☐ Average☐ Breakthrough☐ Connect Four☐ Blind Auction☐ Kuhn Poker☐ Liar's Dice☐ Negotiation☐ Nim☐ Pig☐ Iterated Prisoner's Dilemma☐ Tic-Tac-Toe

Model

☒ GPT-4☒ Llama-2-70b-chat☒ GPT-3.5-turbo☐ Codellama-34b-instruct☐ Mistral-7b-Orca

Agents

☒ ToT☒ Prompt☒ CoT☐ SC-CoT

Opponent Model

☒ GPT-3.5-turbo-1106☐ GPT-4

Opponent Agents

☒ prompt agent

Model	Agent	Opponent Model	Opponent Agent	Average	Breakthrough	Connect Four	Blind Auction	Kuhn Poker	Liar's Dice	Negotiation
GPT-4	Prompt	GPT-3.5-turbo-1106	prompt agent	0.129	0.32	0.08	0.04	0.4	0.8	-0.281
GPT-4	CoT	GPT-3.5-turbo-1106	prompt agent	0.126	0.56	-0.08	0.069	0.44	0.24	0.135
GPT-3.5-turbo	CoT	GPT-3.5-turbo-1106	prompt agent	0.023	-0.12	-0.32	0.115	0.12	0.44	-0.207
GPT-3.5-turbo	Prompt	GPT-3.5-turbo-1106	prompt agent	0	0	0	0	0	0	0
GPT-3.5-turbo	ToT	GPT-3.5-turbo-1106	prompt agent	-0.005	0.16	0.24	-0.12	0	0	0.183
Llama-2-70b-chat	CoT	GPT-3.5-turbo-1106	prompt agent	-0.089	-0.8	0.08	-0.086	-0.2	0.265	0.128
Llama-2-70b-chat	Prompt	GPT-3.5-turbo-1106	prompt agent	-0.205	-0.44	-1	-0.075	-0.04	-0.16	-0.033

Figure A9: The user interface of GTBENCH leaderboard when various LLMs/agents and opponents are selected.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: We summarize our main empirical observations and conclusions in the abstract and introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: Section Limitation

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[NA\]](#)

Justification: This paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We provide details about models and generative hyper-parameters in Section 4, and all the prompts utilized in Appendix A5.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide access to the data and code.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Section 3, Section 4, Appendix A5, A6, A7.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Section 4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The results are obtained from endpoint API providers, e.g., OpenAI (Section 4).

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The research conducted in the paper conform with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Section Impact Statements

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [\[Yes\]](#)

Justification: Section Impact Statements

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [\[Yes\]](#)

Justification: Section 4

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.