

---

# GAMEBENCH: Evaluating Strategic Reasoning Abilities of LLM Agents

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large language models (LLMs) have demonstrated remarkable few-shot performance on many natural language understanding tasks. Despite several demonstrations of using large language models in complex, strategic scenarios, there lacks a comprehensive framework for evaluating agents' performance across various types of reasoning found in games. To address this gap, we introduce GAMEBENCH, a cross-domain benchmark for evaluating strategic reasoning abilities of LLM agents. We focus on 9 different game environments, where each covers at least one axis of key reasoning skill identified in strategy games, and select games for which strategy explanations are unlikely to form a significant portion of models' pretraining corpuses. Our evaluations use GPT-3.5 and GPT-4 in their base form along with two scaffolding frameworks designed to enhance strategic reasoning ability: Chain-of-Thought (CoT) prompting and Reasoning Via Planning (RAP). Our results show that none of the tested models match human performance, and at worse GPT-4 performs worse than random action. CoT and RAP both improve scores but not to comparable human levels. Benchmark code is available at <https://anonymous.4open.science/r/GameBench-5942/>.

## 17 1 Introduction

18 Capabilities of large language models have seen rapid progress, enabling LLMs to be used in agentic tasks (32; 37; 30). This presents opportunities for LLM-based tools to assist humans in several domains, such as API usage (19), web browsing (32) and coding (18). Recent benchmarks have been introduced for evaluating performance on real-world agent tasks (35; 22; 27; 28) with some focused on reasoning (31) and games (21). However, these existing benchmarks are oriented to practical, in-distribution knowledge, which can quickly become saturated with better models.

24 In particular, strategic reasoning is an agentic task that is important for generalising to new contexts, as it involves optimising for an objective in the face of possibly divergent interests of others, where incentives may not be fully known (14). Prior work on reasoning scaffolds also shows that language models have potential to grasp reasoning skills across scenarios (39; 16). Hence, a strategic reasoning benchmark for LLMs, that is inherently multi-agent, would be difficult to saturate. Furthermore, games exemplify environments for demonstrating strategic behaviour in both humans and AI agents, as seen in the well known examples of Chess (34) and Go (33). Hence evaluating LLMs on several types of reasoning behaviours would present a comprehensive, fine-grained benchmark. As such, we introduce GAMEBENCH: a multi-player, cross-domain framework for evaluating strategic reasoning

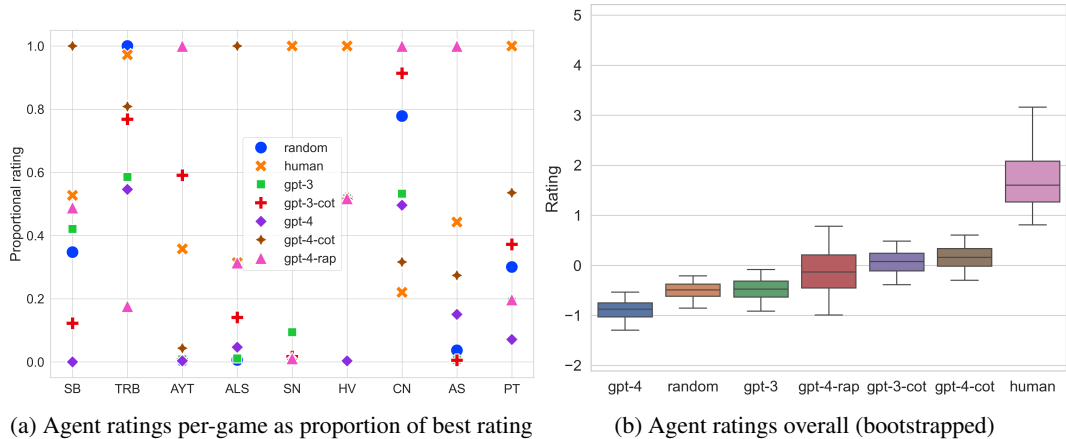


Figure 1: **Rating data** With CoT scaffolding, GPT-4 is the best reasoner below only the human baseline, achieving the best LLM performance on Sea Battle and Pit. But without, it performs worse than even the random baseline due to its exceedingly low rating on Sea Battle. The state-of-the-art RAP scaffolding does not provide as much of an improvement to GPT-4 as CoT does. Looking at the top line of Figure 1a reveals the best agent in each game. See section 3.3 for details. The whiskers represent 90% CIs computed from our bootstrapping process formalized in 3.3. ALS = Air, Land, Sea; ARC = Arctic Scavengers; AYT = Are You the Traitor?; CN = Codenames; HV = Hive; PT = Pit; SN = Santorini; TRB = Two Rooms and a Boom; SB = Sea Battle.

33 in LLM agents using games. We focus on both discrete and open-ended action spaces, across the  
 34 reasoning domains of abstract strategy; non-deterministic outcomes; hidden information; language  
 35 communication; social deduction and cooperation between players. By selecting for games without  
 36 published strategy guides to our knowledge, we ensure that game-specific strategy has been sufficiently  
 37 out-of-distribution in pretraining data. See Table 3 for a complete list of games and game properties.

38 The benchmark consists of obscure board games, card games, and social deception games. We  
 39 evaluate GPT-3.5-turbo-1106 (GPT-3.5) and gpt-4-1106-preview (GPT-4) along with the CoT  
 40 (39) and RAP (16) scaffolding techniques, by playing them against each other, a random-action-  
 41 selector baseline, and a human baseline. We conducted a literature review and identified RAP to be  
 42 the state-of-the-art scaffolding that fit the parameters of our benchmark, i.e. each agent has access to  
 43 the same game state information and no agent can peek at future states. Agents are rated using the  
 44 exponential Bradley-Terry model (6). This has useful advantages over the typical Elo system (13),  
 45 such as its assumption that each agent’s ability is fixed and will not change between matches.

46 Our results show that CoT-augmented and RAP-augmented models demonstrate superior strategic  
 47 superior to the random baseline; that GPT-3.5 matches the random baseline; that GPT-4 performs  
 48 worse than the random baseline; and that the human baseline performs superior to all.

49 With this benchmark, we propose a means to measure the strategic reasoning abilities of LLM agents  
 50 in diverse game environments. Our contributions are as follows:

- 51 • **GAMEBENCH**, the first benchmark to capture both cross-domain and out-of-distribution  
 52 strategic reasoning for comparison across multiple agents.
- 53 • **Empirical results** on GPT-3.5 and GPT-4, demonstrating the effects of Chain-of-Thought  
 54 scaffolding and the state-of-the-art scaffolding.

## 55 2 Related works

56 We provide a detailed literature review of LLM agents playing games, game-theoretic benchmarks,  
 57 dialogue-based benchmarks, and multi-agent game suites in Appendix A.

## 58 3 GAMEBENCH

59 In Section 3.1 we discuss our reasoning behind our selection of agents and scaffolds. In Section  
60 3.2 we describe the agent and game interfaces. In Section 3.3 we introduce our rating model and  
61 formalize our process for calculating ratings. See Appendix I for additional details on how we  
62 selected games to include in the benchmark.

### 63 3.1 Agent and scaffolding selection

64 We benchmark GPT-3.5 (gpt-3.5-turbo-1106) and GPT-4 (gpt-4-1106-preview) due to their  
65 size, mainstream popularity, and convenient public API. We include these base models as well as  
66 several black-box scaffolding interventions in order to measure the relative effects these scaffolding  
67 interventions have on improving the reasoning abilities of the base models. We selected Chain-of-  
68 Thought (39) prompting for its ubiquity and Reasoning-via-Planning (16) for its state-of-the-art status.  
69 We also include a random-action-selecting agent as baseline of no strategic reasoning ability, and a  
70 human agent as a baseline of progress towards human-level strategic reasoning.

71 For more details about agent implementation, see Appendix G.

### 72 3.2 API

73 Each environment, implemented in Python, describes a Game object with methods for initializing,  
74 retrieving the game’s current state and available actions, updating the state with an action, and  
75 executing a full match between two agents. Agents are objects that describe a method for choosing an  
76 action conditioned on the rules, state, and available actions retrieved from a Game instance. Agents are  
77 instantiated at the beginning of a match and destroyed at the end, so agents may maintain persistent  
78 state between moves to choose an action.

### 79 3.3 Rating calculation

80 We formalize our rating calculation as follows. Let our dataset contain  $P$ , the population of all  
81 possible matches across all games, and  $S = \{m_1, m_2, \dots, m_n\}$ , our sample of  $n$  matches. Define  
82 the weight  $w_i$  for each match  $m_i$  to be inversely proportional to the number of matches collected for  
83 that match’s game. Specifically, if match  $m_i$  belongs to game  $X$  which has  $N_X$  matches, then the  
84 weight  $w_i$  is given by:

$$w_i = \frac{1}{N_X}. \quad (1)$$

85 We then perform bootstrapping on the sample  $S$  for  $B = 10,000$  times. Let  $S_b^* =$   
86  $[m_{i_1}, m_{i_2}, \dots, m_{i_n}]$  be the  $b$ th bootstrapped sample, where  $m_{i_j}$  is randomly selected from  $S$  with  
87 probability proportional to  $w_i$  with replacement.

$$P(i > j) = \frac{e^{\beta_i}}{e^{\beta_i} + e^{\beta_j}} \quad (2)$$

88 For each bootstrapped sample  $S_b^*$ , we use maximum-likelihood estimation to fit the parameters of the  
89 above exponential Bradley–Terry model  $\theta_b = \{\beta_{\text{random}}, \beta_{\text{GPT-3.5}}, \dots\}$ . Let  $\theta_{b,k}$  denote the parameter  
90 for agent  $k$  in bootstrapped sample  $b$ . We take the means of these distributions to be the “true” rating  
91  $\hat{\theta}_k$  for each agent  $k$ , given by:

$$\hat{\theta}_k = \frac{1}{B} \sum_{b=1}^B \theta_{b,k} \quad (3)$$

92 We considered several methods for aggregating pairwise match results across games into scores  
93 that represent the general skill of each model, including the Elo system (13). Unlike Elo, the  
94 Bradley–Terry system (6) assumes model skill does not change over time and does not need to be  
95 calculated in a decentralized manner, making it more suitable for evaluating language models (10).

96 **4 Empirical results**

97 Additional figures showing agent-pairwise data covering the number of games, total score, win  
 98 probability, and rating per game is available in Appendix L. The rating plots in Appendix L show  
 99 90% confidence intervals for the points in Figure 1a.

Table 1: **Game ratings** The table highlights the effects of scaffolds. Across all games, GPT-4 with CoT scaffolding improves over the base model substantially. But GPT-3.5 with CoT scaffolding is outperformed by the base model in *ALS*, *HV*, and *TRB*. Additionally, GPT-4 with RAP scaffolding usually under-performs GPT-4-CoT except in *AYT*, *SB*, and *TRB*.

Agent	Rating									
	Overall	ALS	ARC	AYT	CN	HV	PT	SN	TRB	SB
Random	-0.50	1.07	<b>0.48</b>	-2.52	-2.67	-1.15	<u>0.63</u>	0.37	-0.79	0.05
Human	<b>1.76</b>	<u>1.49</u>	<u>0.45</u>	1.92	1.26	<b>3.63</b>	<b>1.29</b>	-0.89	<u>1.70</u>	<b>1.25</b>
GPT-3.5	-0.48	<u>1.26</u>	-0.05	-1.84	-2.06	<u>1.27</u>	<u>0.63</u>	-0.01	-2.51	-0.41
GPT-3.5-CoT	0.06	0.03	0.22	<u>2.42</u>	0.45	-0.44	<u>0.63</u>	<u>0.53</u>	-2.76	0.26
GPT-4	-0.89	-7.38	-0.12	-2.73	-0.65	-1.31	-4.42	-0.08	0.62	-1.40
GPT-4-CoT	<u>0.16</u>	<b>2.13</b>	0.27	-0.19	<b>2.41</b>	-1.13	0.63	-0.53	1.22	<u>0.62</u>
GPT-4-RAP	-0.10	1.41	-1.25	<b>2.94</b>	<u>1.26</u>	-0.86	<u>0.63</u>	<b>0.62</b>	<b>2.51</b>	-0.37

Table 2: **Average score.** The total score an agent achieved in a game divided by the number of games that agent played. Comparing with 1, this table highlights interesting correlations between empirical score and model-inferred ratings. For example, in *Air*, *Land*, and *Sea*, GPT-4-CoT has the top rating while the human baseline has second-top, but they swap when examining average score. This plot also shows more clearly why the human baseline has the highest rating even though both the human baseline and GPT-4-RAP have the highest rating in three games. Here, the human baseline achieved the highest score in four games but GPT-4-RAP only achieved the highest in two.

Agent	Score									
	Overall	ALS	ARC	AYT	CN	HV	PT	SN	TRB	SB
Random	0.49	0.72	<b>0.60</b>	0.25	0.18	0.41	<u>0.50</u>	0.56	0.52	<u>0.58</u>
Human	<b>0.85</b>	<b>1.00</b>	NaN	NaN	NaN	<b>1.00</b>	<b>1.00</b>	0.43	NaN	<b>0.78</b>
GPT-3.5	0.48	0.64	0.43	0.43	0.63	<u>0.80</u>	<u>0.50</u>	0.47	0.27	0.40
GPT-3.5-COT	0.60	0.43	<u>0.50</u>	<u>0.93</u>	<u>0.89</u>	0.60	<u>0.50</u>	<b>0.61</b>	0.33	0.55
GPT-4	0.31	0.00	0.42	0.33	0.83	0.33	0.31	0.42	0.71	0.20
GPT-4-COT	0.60	<u>0.81</u>	<u>0.50</u>	0.64	<b>1.00</b>	0.50	<u>0.50</u>	0.37	<u>0.75</u>	0.51
GPT-4-RAP	<u>0.62</u>	NaN	0.33	<b>1.00</b>	NaN	0.50	NaN	<u>0.58</u>	<b>1.00</b>	0.26

100 **5 Conclusion**

101 We present GAMEBENCH, an LLM agent benchmark to test strategic reasoning ability using diverse  
 102 games that have sparse strategy material in pretraining data. We benchmark OpenAI’s GPT-3.5 and  
 103 GPT-4 models and evaluate the impact of two scaffolding methods: Chain of Thought (CoT) and  
 104 Reasoning via Planning (RAP). We find that human trials consistently outperform all LLM agents.  
 105 Of all the agent configurations, CoT agents performed the best, followed by RAP-augmented GPT-4.  
 106 Base GPT-3.5 performed on-par with the random baseline, and base GPT-4 performed worse. These  
 107 results show that while measures such as scaffolding can help improve performance in strategic  
 108 reasoning, even the best configuration fall short of human reasoning. LLMs show great promise  
 109 working on in-distribution tasks, though their performance on OOD task sets show a low risk for  
 110 current dangers of deploying autonomous agents. Nonetheless, the performance gains achieved  
 111 through scaffolding techniques indicate the potential for future advancements that could increase the  
 112 risk posed by such systems if their reasoning capabilities continue to improve.

113 **Acknowledgments and Disclosure of Funding**

114 We thank Joshua Clymer for providing advisory help. We thank Severin Field and Misha Gerovitch  
115 for providing feedback on our drafts. We thank Shubhorup Biswas for implementing Atari Boxing.

116 **References**

- 117 [1] ABDELNABI, S., GOMAA, A., SIVAPRASAD, S., SCHONHERR, L., AND FRITZ, M.  
118 Llm-deliberation: Evaluating llms with interactive multi-agent negotiation games. ArXiv  
119 abs/2309.17234 (2023).
- 120 [2] ABDULHAI, M., WHITE, I., SNELL, C. B., SUN, C., HONG, J., ZHAI, Y., XU, K., AND  
121 LEVINE, S. Lmrl gym: Benchmarks for multi-turn reinforcement learning with language  
122 models. ArXiv abs/2311.18232 (2023).
- 123 [3] AGASHE, S., FAN, Y., REYNA, A., AND WANG, X. E. Llm-coordination: Evaluating and  
124 analyzing multi-agent coordination abilities in large language models, 2024.
- 125 [4] AKATA, E., SCHULZ, L., CODA-FORNO, J., OH, S. J., BETHGE, M., AND SCHULZ, E.  
126 Playing repeated games with large language models. ArXiv abs/2305.16867 (2023).
- 127 [5] BAKHTIN, A., BROWN, N., DINAN, E., FARINA, G., FLAHERTY, C., FRIED, D., GOFF,  
128 A., GRAY, J., HU, H., JACOB, A. P., KOMEILI, M., KONATH, K., KWON, M., LERER, A.,  
129 LEWIS, M., MILLER, A. H., MITTS, S., RENDUCHINTALA, A., ROLLER, S., ROWE, D.,  
130 SHI, W., SPISAK, J., WEI, A., WU, D. J., ZHANG, H., AND ZIJLSTRA, M. Human-level play  
131 in the game of diplomacy by combining language models with strategic reasoning. Science **378**  
132 (2022), 1067 – 1074.
- 133 [6] BRADLEY, R. A., AND TERRY, M. E. Rank analysis of incomplete block designs: I. the  
134 method of paired comparisons. Biometrika **39**, 3/4 (1952), 324–345.
- 135 [7] CHALAMALASETTI, K., GOTZE, J., HAKIMOV, S., MADUREIRA, B., SADLER, P., AND  
136 SCHLANGEN, D. clembench: Using game play to evaluate chat-optimized language models as  
137 conversational agents. ArXiv abs/2305.13455 (2023).
- 138 [8] CHEN, J., HU, X., LIU, S., HUANG, S., TU, W., HE, Z., AND WEN, L. Llmarena: As-  
139 ssuming capabilities of large language models in dynamic multi-agent environments. ArXiv  
140 abs/2402.16499 (2024).
- 141 [9] CHEN, J., YUAN, S., YE, R., MAJUMDER, B. P., AND RICHARDSON, K. Put your money  
142 where your mouth is: Evaluating strategic planning and execution of llm agents in an auction  
143 arena. ArXiv abs/2310.05746 (2023).
- 144 [10] CHIANG, W.-L., ZHENG, L., SHENG, Y., ANGELOPOULOS, A. N., LI, T., LI, D., ZHANG,  
145 H., ZHU, B., JORDAN, M., GONZALEZ, J. E., AND STOICA, I. Chatbot arena: An open  
146 platform for evaluating llms by human preference. arXiv preprint arXiv:2403.04132 (2023).
- 147 [11] CHOWDHERY, A., NARANG, S., DEVLIN, J., BOSMA, M., MISHRA, G., ROBERTS,  
148 A., BARHAM, P., CHUNG, H. W., SUTTON, C., GEHRMANN, S., SCHUH, P., SHI, K.,  
149 TSVYASHCHENKO, S., MAYNEZ, J., RAO, A., BARNES, P., TAY, Y., SHAZEER, N., PRAB-  
150 HAKARAN, V., REIF, E., DU, N., HUTCHINSON, B., POPE, R., BRADBURY, J., AUSTIN,  
151 J., ISARD, M., GUR-ARI, G., YIN, P., DUKE, T., LEVSKAYA, A., GHEMAWAT, S., DEV,  
152 S., MICHALEWSKI, H., GARCIA, X., MISRA, V., ROBINSON, K., FEDUS, L., ZHOU, D.,  
153 IPPOLITO, D., LUAN, D., LIM, H., ZOPH, B., SPIRIDONOV, A., SEPASSI, R., DOHAN, D.,  
154 AGRAWAL, S., OMERNICK, M., DAI, A. M., PILLAI, T. S., PELLAT, M., LEWKOWYCZ,  
155 A., MOREIRA, E., CHILD, R., POLOZOV, O., LEE, K., ZHOU, Z., WANG, X., SAETA, B.,  
156 DIAZ, M., FIRAT, O., CATASTA, M., WEI, J., MEIER-HELLSTERN, K., ECK, D., DEAN, J.,  
157 PETROV, S., AND FIEDEL, N. Palm: Scaling language modeling with pathways, 2022.

- 158 [12] DUAN, J., ZHANG, R., DIFFENDERFER, J., KAILKHURA, B., SUN, L., STENGEL-ESKIN, E.,  
159 BANSAL, M., CHEN, T., AND XU, K. Gtbench: Uncovering the strategic reasoning limitations  
160 of llms via game-theoretic evaluations. [ArXiv abs/2402.12348](https://arxiv.org/abs/2402.12348) (2024).
- 161 [13] ELO, A. E. The proposed uscf rating system, its development, theory, and applications. [Chess  
162 life 22](https://doi.org/10.1145/254181.254182), 8 (1967), 242–247.
- 163 [14] GANDHI, K., SADIGH, D., AND GOODMAN, N. D. Strategic reasoning with language models,  
164 2023.
- 165 [15] GANDHI, K., SADIGH, D., AND GOODMAN, N. D. Strategic reasoning with language models.  
166 [ArXiv abs/2305.19165](https://arxiv.org/abs/2305.19165) (2023).
- 167 [16] HAO, S., GU, Y., MA, H., HONG, J. J., WANG, Z., WANG, D. Z., AND HU, Z. Reasoning  
168 with language model is planning with world model. [arXiv preprint arXiv:2305.14992](https://arxiv.org/abs/2305.14992) (2023).
- 169 [17] ILIĆ, D. Unveiling the general intelligence factor in language models: A psychometric approach,  
170 2023.
- 171 [18] KAZEMITABAAR, M., HOU, X., HENLEY, A., ERICSON, B. J., WEINTROP, D., AND GROSS-  
172 MAN, T. How novices use llm-based code generators to solve cs1 coding tasks in a self-paced  
173 learning environment, 2023.
- 174 [19] LI, M., ZHAO, Y., YU, B., SONG, F., LI, H., YU, H., LI, Z., HUANG, F., AND LI, Y.  
175 Api-bank: A comprehensive benchmark for tool-augmented llms, 2023.
- 176 [20] LIGHT, J., CAI, M., SHEN, S., AND HU, Z. Avalonbench: Evaluating llms playing the game  
177 of avalon, 2023.
- 178 [21] LIN, J., ZHAO, H., ZHANG, A., WU, Y., PING, H., AND CHEN, Q. Agentsims: An open-  
179 source sandbox for large language model evaluation, 2023.
- 180 [22] LIU, X., YU, H., ZHANG, H., XU, Y., LEI, X., LAI, H., GU, Y., DING, H., MEN, K., YANG,  
181 K., ZHANG, S., DENG, X., ZENG, A., DU, Z., ZHANG, C., SHEN, S., ZHANG, T., SU, Y.,  
182 SUN, H., HUANG, M., DONG, Y., AND TANG, J. Agentbench: Evaluating llms as agents,  
183 2023.
- 184 [23] LIU, X., YU, H., ZHANG, H., XU, Y., LEI, X., LAI, H., GU, Y., GU, Y., DING, H., MEN,  
185 K., YANG, K., ZHANG, S., DENG, X., ZENG, A., DU, Z., ZHANG, C., SHEN, S., ZHANG,  
186 T., SU, Y., SUN, H., HUANG, M., DONG, Y., AND TANG, J. Agentbench: Evaluating llms as  
187 agents. [ArXiv abs/2308.03688](https://arxiv.org/abs/2308.03688) (2023).
- 188 [24] MAITRIX ORG. Llm-reasoners: A library for advanced large language model reasoning. [https:  
189 //github.com/matrix-org/llm-reasoners](https://github.com/matrix-org/llm-reasoners), 2023. GitHub repository, accessed: 2024-06-  
190 04.
- 191 [25] MAO, S., CAI, Y., XIA, Y., WU, W., WANG, X., WANG, F., GE, T., AND WEI, F. Alympics:  
192 Llm agents meet game theory – exploring strategic decision-making with ai agents, 2023.
- 193 [26] MAYSTRE, L. choix: Inference algorithms for models based on luce’s choice axiom. [https:  
194 //github.com/lucasmaystre/choix](https://github.com/lucasmaystre/choix), 2015. GitHub repository, accessed: 2024-06-04.
- 195 [27] METR. Evaluating language-model agents on realistic autonomous tasks. [https://metr.  
196 org/blog/2023-08-01-new-report/](https://metr.org/blog/2023-08-01-new-report/), 2023.
- 197 [28] MIALON, G., FOURRIER, C., SWIFT, C., WOLF, T., LECUN, Y., AND SCIALOM, T. Gaia: a  
198 benchmark for general ai assistants, 2023.
- 199 [29] QIAO, D., WU, C., LIANG, Y., LI, J., AND DUAN, N. Gameeval: Evaluating llms on  
200 conversational games. [ArXiv abs/2308.10032](https://arxiv.org/abs/2308.10032) (2023).

- 201 [30] RICHARDS, T. B. Autogpt: An autonomous gpt-4 experiment. [https://github.com/](https://github.com/Significant-Gravitas/AutoGPT/tree/master)  
202 Significant-Gravitas/AutoGPT/tree/master, 2023.
- 203 [31] SAWADA, T., PALEKA, D., HAVRILLA, A., TADEPALLI, P., VIDAS, P., KRANIAS, A., NAY,  
204 J. J., GUPTA, K., AND KOMATSUZAKI, A. Arb: Advanced reasoning benchmark for large  
205 language models, 2023.
- 206 [32] SCHICK, T., DWIVEDI-YU, J., DESSI, R., RAILEANU, R., LOMELI, M., ZETTLEMOYER, L.,  
207 CANCEDDA, N., AND SCIALOM, T. Toolformer: Language models can teach themselves to  
208 use tools, 2023.
- 209 [33] SILVER, D., HUANG, A., MADDISON, C., GUEZ, A., SIFRE, L., DRIESSCHE, G., SCHRIT-  
210 TWIESER, J., ANTONOGLU, I., PANNEERSHELVAM, V., LANCTOT, M., DIELEMAN, S.,  
211 GREWE, D., NHAM, J., KALCHBRENNER, N., SUTSKEVER, I., LILICRAP, T., LEACH, M.,  
212 KAVUKCUOGLU, K., GRAEPEL, T., AND HASSABIS, D. Mastering the game of go with deep  
213 neural networks and tree search. *Nature* 529 (01 2016), 484–489.
- 214 [34] SILVER, D., HUBERT, T., SCHRITTWIESER, J., ANTONOGLU, I., LAI, M., GUEZ, A.,  
215 LANCTOT, M., SIFRE, L., KUMARAN, D., GRAEPEL, T., LILICRAP, T., SIMONYAN, K.,  
216 AND HASSABIS, D. Mastering chess and shogi by self-play with a general reinforcement  
217 learning algorithm, 2017.
- 218 [35] WANG, Y., MA, X., ZHANG, G., NI, Y., CHANDRA, A., GUO, S., REN, W., ARULRAJ, A.,  
219 HE, X., JIANG, Z., LI, T., KU, M., WANG, K., ZHUANG, A., FAN, R., YUE, X., AND CHEN,  
220 W. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark,  
221 2024.
- 222 [36] WANG, Z., CAI, S., LIU, A., MA, X., AND LIANG, Y. Describe, explain, plan and select:  
223 Interactive planning with large language models enables open-world multi-task agents. *ArXiv*  
224 [abs/2302.01560](https://arxiv.org/abs/2302.01560) (2023).
- 225 [37] WATKINS, A., SUBEDI, S., AND SHRESTHA, A. Agentgpt. [https://github.com/reworkd/](https://github.com/reworkd/AgentGPT)  
226 AgentGPT, 2023.
- 227 [38] WEI, J., TAY, Y., BOMMASANI, R., RAFFEL, C., ZOPH, B., BORGEAUD, S., YOGATAMA,  
228 D., BOSMA, M., ZHOU, D., METZLER, D., CHI, E. H., HASHIMOTO, T., VINYALS, O.,  
229 LIANG, P., DEAN, J., AND FEDUS, W. Emergent abilities of large language models, 2022.
- 230 [39] WEI, J., WANG, X., SCHUURMANS, D., BOSMA, M., ICHTER, B., XIA, F., CHI, E., LE, Q.,  
231 AND ZHOU, D. Chain-of-thought prompting elicits reasoning in large language models. *arXiv*  
232 [preprint arXiv:2201.11903](https://arxiv.org/abs/2201.11903) (2022).
- 233 [40] WU, Y., TANG, X., MITCHELL, T. M., AND LI, Y. Smartplay : A benchmark for llms as  
234 intelligent agents. *ArXiv abs/2310.01557* (2023).
- 235 [41] WU, Z., QIU, L., ROSS, A., AKYÜREK, E., CHEN, B., WANG, B., KIM, N., ANDREAS, J.,  
236 AND KIM, Y. Reasoning or reciting? exploring the capabilities and limitations of language  
237 models through counterfactual tasks. *ArXiv abs/2307.02477* (2023).
- 238 [42] XU, L., HU, Z., ZHOU, D., REN, H., DONG, Z., KEUTZER, K., NG, S.-K., AND FENG, J.  
239 Magic: Investigation of large language model powered multi-agent in cognition, adaptability,  
240 rationality and collaboration. *ArXiv abs/2311.08562* (2023).
- 241 [43] XU, Y., WANG, S., LI, P., LUO, F., WANG, X., LIU, W., AND LIU, Y. Exploring  
242 large language models for communication games: An empirical study on werewolf. *ArXiv*  
243 [abs/2309.04658](https://arxiv.org/abs/2309.04658) (2023).
- 244 [44] ZELIKMAN, E., WU, Y., MU, J., AND GOODMAN, N. D. Star: Bootstrapping reasoning with  
245 reasoning, 2022.

246 [45] ZHU, X., CHEN, Y., TIAN, H., TAO, C., SU, W., YANG, C., HUANG, G., LI, B., LU, L.,  
247 WANG, X., QIAO, Y., ZHANG, Z., AND DAI, J. Ghost in the minecraft: Generally capable  
248 agents for open-world environments via large language models with text-based knowledge and  
249 memory. [ArXiv abs/2305.17144](https://arxiv.org/abs/2305.17144) (2023).

## 250 Checklist

- 251 1. For all authors...
  - 252 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
253 contributions and scope? [Yes] The paper's contributions are listed plainly at the end  
254 of the introduction and explained in-depth in sections 3 and 4.
  - 255 (b) Did you describe the limitations of your work? [Yes] See Section B.
  - 256 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See  
257 Appendix D.
  - 258 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
259 them? [Yes] We discuss human research in Appendix E, data concerns in Appendix F,  
260 and societal impact in Appendix D.
- 261 2. If you are including theoretical results...
  - 262 (a) Did you state the full set of assumptions of all theoretical results? [N/A] We are not  
263 including theoretical results.
  - 264 (b) Did you include complete proofs of all theoretical results? [N/A] We are not including  
265 theoretical results.
- 266 3. If you ran experiments (e.g. for benchmarks)...
  - 267 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
268 mental results (either in the supplemental material or as a URL)? [Yes] See Appendix  
269 G.
  - 270 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
271 were chosen)? [N/A] We did not do any training.
  - 272 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
273 ments multiple times)? [Yes] Figure 1b shows error bars, which are explained in the  
274 caption. And Appendix L shows more rating breakdowns by game with error bars.
  - 275 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
276 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix F.
- 277 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - 278 (a) If your work uses existing assets, did you cite the creators? [Yes] We cite creators in  
279 our Github and in Appendix G.
  - 280 (b) Did you mention the license of the assets? [Yes] We mention licenses in our Github  
281 where required and in Appendix G.
  - 282 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
283 See Appendix F.
  - 284 (d) Did you discuss whether and how consent was obtained from people whose data you're  
285 using/curating? [N/A] We are not using nor curating other people's data.
  - 286 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
287 information or offensive content? [Yes] See Appendix E.
- 288 5. If you used crowdsourcing or conducted research with human subjects...
  - 289 (a) Did you include the full text of instructions given to participants and screenshots, if  
290 applicable? [Yes] See Appendix E.
  - 291 (b) Did you describe any potential participant risks, with links to Institutional Review  
292 Board (IRB) approvals, if applicable? [Yes] See Appendix E.



293 (c) Did you include the estimated hourly wage paid to participants and the total amount  
294 spent on participant compensation? [Yes] See Appendix E.

## 295 A Related works

296 **LLM agents playing games** Using games to evaluate LLMs has significant precedent in previous  
297 research. Some studies evaluate models using single strategic tasks or games, such as Minecraft  
298 (36; 45), Diplomacy (5), Avalon (20), and Werewolf (43). Other benchmarks (40; 23) capture a more  
299 comprehensive picture by using suites of multiple tasks or games to evaluate LLMs as intelligent  
300 agents. However, the tasks represented in these benchmarks don't involve interaction with other  
301 agents, so they don't reflect strategic reasoning as defined in this work.

302 **Game-theoretic scenarios** Several benchmark suites focus on common game theory scenarios, such  
303 as auctions (9; 25), matrix games like Prisoner's Dilemma (4; 15), and negotiation (1; 15). While  
304 they do involve multi-agent interaction and are useful for testing models' strategic reasoning ability,  
305 our benchmark focuses on more complex games that aren't as frequently studied as these game theory  
306 scenarios. Given no major strategy guides or forums dedicated to these games, we believe there is  
307 less documentation on optimal strategies for them present in LLMs' training corpuses.

308 **Dialogue-based games** Some benchmarks employ dialogue-based games that are less well-  
309 documented on the internet: (**author?**) (3) and (**author?**) (7) use novel cooperative dialogue games,  
310 and (**author?**) (29) uses two social deduction games and one word guessing game. However, our  
311 benchmark aims to evaluate LLMs' strategic reasoning ability not only in cooperative and conversa-  
312 tional environments, but competitive, spatial, and non-deterministic ones as well.

313 **Diverse multi-agent game suites** The benchmarks most similar to ours employ diverse suites of  
314 complex multi-agent games, including conversational, board, and card games (8; 12; 2; 42). However,  
315 many of the included games are either commonly found on the internet, such as TicTacToe, Poker,  
316 and Connect Four, or common game-theoretic scenarios, as discussed previously. These games are  
317 not as out-of-distribution as desired.

318 In summary, we build upon previous work by introducing a diverse suite of multi-agent games  
319 to evaluate the strategic reasoning ability of LLMs as agents. Our benchmark is characterized by  
320 its inclusion of complex games that span a range of game characteristics and are not likely to be  
321 well-represented in LLMs' pretraining corpuses.

## 322 B Discussion

323 We now discuss the limitations and future directions of our work.

324 **Confirming out-of-distribution status** It is clear by simply asking GPT-4 that it already knows  
325 about these games and their rules. It is unclear, however, if it consumed any strategy guides about  
326 these games in the pretraining process, which is the determining factor for out-of-distribution status  
327 in our benchmark. **Future work** We propose the following experiment. Design an intervention that  
328 is: supply a strategy guide in-context to a language model agent for the game it is playing. We  
329 would expect this intervention to improve agent performance more on out-of-distribution games  
330 than in-distribution games. Collect data of agents playing an unknown distribution game; agents  
331 with the intervention playing an unknown distribution game; agents playing known in-distribution  
332 games; agents with the intervention playing known in-distribution games. Compare the effect of  
333 the intervention on the unknown distribution game versus the effect on the known in-distribution  
334 games. If the effect is much higher on the unknown distribution game, this is a evidence for the game  
335 being out-of-distribution. This would work better with known out-of-distribution games, but this  
336 may not be possible to know in all cases. We could also compare models' performance on common  
337 games vs. "counterfactual" games, which are slightly modified to reduce any association with their  
338 in-distribution counterparts (41).

339 **Protecting out-of-distribution status** We did not attempt to protect these games from becoming  
340 in-distribution in the future. **Future work** Developers of frontier models should curate strategic  
341 reasoning environments by ensuring these games are held out from pretraining data. For ubiquitous  
342 games such as chess, this is unfeasible. But following our heuristics for game selection discussed in  
343 Appendix I, it should be reasonable to find games without much internet data.

344 **Results’ sensitivity to games** From inspecting GPT-4’s surprisingly low rating with *Sea Battle*, it  
345 became apparent that our “multigame” approach to aggregation may be inadequate due its sensitivity  
346 to the games included; i.e., ablating *Sea Battle* significantly changed the data narrative. **Future**  
347 **work** We see multiple ways forward. If aggregate data is useful, investigate more robust forms  
348 of aggregation, such as the g-factor or factor analysis in general. Alternatively, explore a multi-  
349 dimensional approach that attempts to score agents on the six reasoning categories from Table 3.  
350 Or, discard any notion of aggregation and determine effective means of analysis that looks only at  
351 individual games and maybe uses more qualitative data with the help of human experts.

352 **Low-resolution human benchmark** We find it especially important to know how well these models  
353 fair compared to humans, but collecting comprehensive human data was out of our means. **Future**  
354 **work** Conduct more comprehensive human data to form a distribution of human strengths on each  
355 game with which we can measure the progress of model and scaffolding development.

356 **Uncaught edge cases** Every few games were inspected during data collection, and occasionally, we  
357 caught and fixed bugs in our evaluation code. It is possible that some edge cases went unnoticed and  
358 were featured in our final data release. **Future work** Incorporating more human subjects into the data  
359 collection should make this process trivial, as they can provide immediate feedback if they witness  
360 unexpected behavior.

361 **Benchmark and dataset size** Our benchmark has a respectable number of games and agents  
362 compared to other benchmarks (8; 12; 2; 42), but the addition of more games and agents would  
363 provide a richer picture of models’ strategic reasoning abilities. Additionally, our dataset is fairly  
364 small and suffers from biases from variable resource cost between games. **Future work** Add more  
365 varied games to the benchmark, evaluate more model and scaffolding configurations, and collect  
366 more data for each configuration.

## 367 C Interpreting empirical results

### 368 C.1 Human comparison

369 The human baseline outperforms all model and scaffolding configurations in the benchmark. The  
370 upper-bound of GPT-4-RAP’s confidence interval in Figure 1b just reaches the lower-bound of the  
371 human baseline. But due to both GPT-4-RAP and the human baseline having very few data points,  
372 this detail should not be taken very seriously. In Table 2, the human baseline achieves the highest  
373 overall score in every game it played except for *Santorini*.

374 The human subject beat their opponent agent in all matches except for two of the three *Codenames*  
375 matches. For these particular matches, the human subject employed a friend because *Codenames*  
376 typically requires at least two players per team. We hypothesize that LLM agents perform better  
377 in this context because they are better at modeling their teammate’s thought process, as they are  
378 instantiated from the same underlying language model. In contrast, pairs of humans share much less  
379 cognitive similarity.

380 Details about the human data collection process are discussed in Appendix E.

### 381 C.2 Effect of scaffolding

382 Chain-of-Thought prompting provided the best median and upper quartile results of all configurations  
383 tested in Figure 1b. GPT-3.5 and GPT-4 showed almost identical performance with GPT-4 with

384 only a slight improvement over GPT-3.5. The positive effects of CoT prompting are already well-  
385 documented (11; 44; 38), and our results provides evidence of their use in strategic settings.

386 If we interpret the addition of CoT scaffolding as an intervention on the base model, we see it  
387 improves strategic reasoning ability in GPT-4 moreso than in GPT-3.5. In *Sea Battle*, this intervention  
388 brings GPT-4 from the worst model to the best model. In every game except *Codenames*, GPT-4 with  
389 CoT scaffolding outperforms its base model. But for GPT-3.5, the base model outperforms the CoT  
390 variant in *Santorini* and *Sea Battle*. One possible hypothesis for this difference in effect between  
391 on GPT-3.5 and GPT-4 is that GPT-4 is a bigger model and thus can probably make better use of  
392 in-context information.

### 393 **C.3 GPT-3.5 versus GPT-4**

394 GPT-3.5 performs only slightly better than random action. Surprisingly, GPT-4 performs the worst  
395 of all configurations with its upper quartile performance being worse than random’s lowest quartile.  
396 This result is mostly due to GPT-4 losing all matches in *Sea Battle*. This challenges our aggregation  
397 method: GPT-4 should not be so harshly penalized for poor performance on one game.

398 An alternative aggregation method that would be more robust to outliers is to use factor analysis  
399 to isolate a “general strategic reasoning factor” that explains a significant portion of the variance  
400 between models’ performances. This method is used to aggregate separate cognitive test scores  
401 into IQ scores, making it apt for evaluating LLMs’ reasoning abilities (17). We expect this g-factor  
402 approach to appropriately weigh models’ *Sea Battle* ratings lower, fixing this discrepancy.

403 Considering these two datapoints and analysis from C.2, we can tentatively conclude that strategic  
404 reasoning ability is not improving in OpenAI’s newest frontier models alone, but their receptiveness  
405 to scaffolding to improve strategic reasoning is increasing.

### 406 **C.4 State-of-the-art scaffolding**

407 The state-of-the-art scaffolding was outperformed by both Chain-of-Thought agents. One possible  
408 hypothesis for this is that, during the Monte-Carlo tree search, this agent predicts new states based  
409 on the state being examined, which is already a predicted state depending if  $\text{depth} \geq 1$ . If the  
410 agent makes any errors in this examined state’s prediction due to misunderstandings about the game  
411 state or rules, these will likely be compounded in the next set of predictions. We might expect the  
412 Chain-of-Thought agents to be susceptible to the same issue of compounding errors, but to a lesser  
413 extent. This could be tested qualitatively by a human expert analysing GPT-4-RAP’s predictions for  
414 accuracy.

415 Another hypothesis is that we ran GPT-4-RAP to a depth great-enough to surpass GPT-4 without  
416 RAP scaffolding, but not great-enough depth to surpass Chain-of-Thought scaffolding. This could be  
417 tested by adding several GPT-4-RAP agents to the benchmark, each with different depths.

418 It seems unlikely that Chain-of-Thought prompting should be the most sophisticated black-box  
419 scaffolding, so it remains an open question to find this scaffolding in order to establish an upper-  
420 bound on strategic reasoning ability with black-box scaffolds.

## 421 **D Hazards**

422 We believe that good strategic reasoning is a dangerous capability for an AI agent to have, especially  
423 one that will operate autonomously. Thus, good performance on this benchmark could correlate with  
424 harmful risk. This is important for labs developing frontier models to be able to measure and be  
425 aware of, but it is also possible for a malicious or ignorant actor to use this benchmark as a feedback  
426 signal to improve their own large language model’s strategic reasoning ability. However, we think  
427 that the former benefit outweighs the latter risk in this time where the development of large language  
428 models is largely controlled by a few frontier labs. And we reduce the risk from ignorant actors by  
429 producing these benchmarks and discussing their importance.

## 430 **E Research on human subjects**

431 Our human-based data-points came from a co-creator of the benchmark, and the same person with  
432 their friend for *Codenames*.

433 The instructions were communicated informally because the subject co-designed the benchmark and  
434 this human study. They were initially instructed to play against the GPT-4-RAP, but due to resource  
435 costs, were later instructed to play against the any agent except GPT-4-RAP or the random baseline.  
436 They were instructed to not play *Are You the Traitor?* and *Two Rooms and a Boom* because they are  
437 social deduction games and it is not a good setup to have one agent with extra information than other  
438 agents. They were instructed to collect as many matches as they were willing to collect in the time  
439 they had available.

440 No additional compensation was provided to them for data collection, but the API costs were covered.

441 Given the informal nature of the data collection, the near-zero risk, and the fact that the subject was a  
442 co-creator in this benchmark and this experiment, we did not discuss risks or consult an IRB. The  
443 data this person created do not contain any identifying information.

## 444 **F Dataset documentation**

445 The data used to generate the figures and tables in this paper are available in our Github <https://anonymous.4open.science/r/GameBench-5942/> under the CC-BY 4.0 license. These data will  
446 remain available here as long as Github is available. New data may be added by the authors in the  
447 future, which will be documented in the commits.

449 The intended use of this data is to compare GPT-3.5 and GPT-4 on this benchmark, and to compare  
450 against new models, scaffolds, baselines, and informed-and-consenting humans in the future.

451 The data are in JSON format. The top-level object is an array, and the array contains objects. Each  
452 object has a “game” key which indicates the game, and two other keys – the two agents that played in  
453 no particular order – with their respective score as the value. Scores are in the range  $[0, 1]$  and sum  
454 to 1.

455 Our data collection was not uniform across games nor against agent-pairs due to resource constraints.  
456 In general, we preferred playing agents against the random baseline and preferred games that didn’t  
457 take too long to complete.

458 All data for each agent except the random agent were collected using OpenAI’s completions API.  
459 Each game was designed to take no more than 5 minutes when playing base GPT-4 against random.  
460 Cost estimates were not obtained, but it can be assumed that CoT agents will cost approximately  
461 twice their base variants, and GPT-4-RAP will cost approximately base cost  $\times$  MCTS depth  $\times$   
462 number of actions per state  $\times$  6

## 463 **G Additional implementation details**

464 To measure multimodal capabilities, *Hive* was made to use images to represent its game state, instead  
465 of text like all the other games. However, GPT-3.5 is not multimodal, so it was served textual  
466 representations of the graphical state created by GPT-4. Then, for RAP, GPT-4 with the completions  
467 API can’t produce images when predicting future states, so for simplicity, the image is turned into a  
468 text description here as well.

469 GPT-4-RAP was run with the default parameters from the `llm-reasoners` library (24) except the  
470 Monte-Carlo tree search depth limit was set to 2 due to resource constraints.

471 Data was not collected for a GPT-3.5-RAP because GPT-3.5 refused to comply with prompts asking  
472 it to predict actions, game states, or other players’ behaviors. The model would often reply, “As a  
473 language model, I can incapable of predicting...” Because it is unlikely that GPT-3.5 is self-aware and

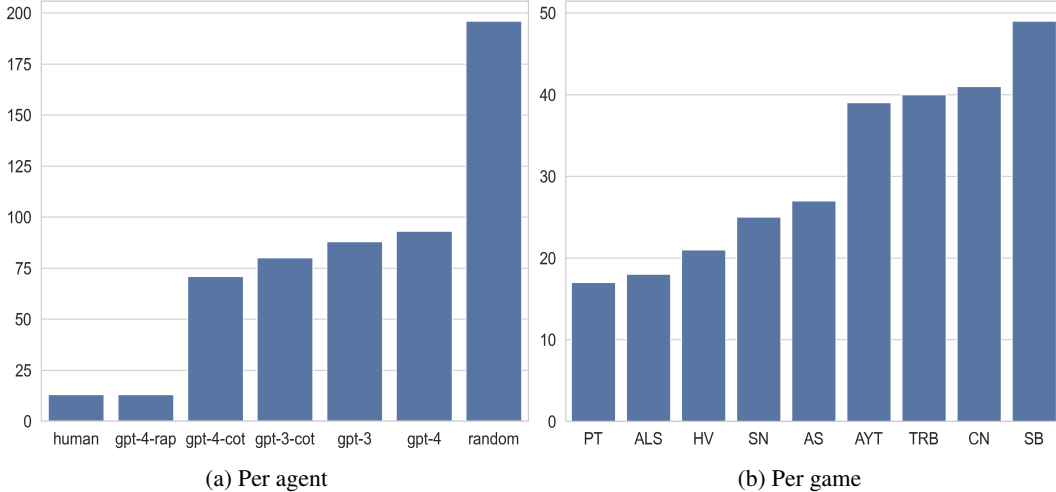


Figure 2: **Number of matches recorded** The random baseline and faster games were oversampled due to their low cost.

474 because we never indicate to the model that it is a language model in our prompting, we hypothesize  
 475 that this refusal is due to the nature of GPT-3.5’s hidden system prompting for refusing unsafe  
 476 behaviors and not due to a lack of ability on GPT-3.5’s part. As such, it is difficult to measure the  
 477 relative effect of RAP scaffolding on GPT-4 versus GPT-3.5.

478 CoT-scaffolded agents are prompted with “First, let’s reason out loud about which action  
 479 you should take to maximize your probability of winning.” after they see the game state  
 480 and available actions.

481 GPT-4-RAP employs a Monte-Carlo tree search where states and actions are model predictions, and  
 482 rewards are computed using next-token probabilities. Our implementation of RAP relies heavily on  
 483 code from (16). Their code is available under the Apache License 2.0. Details of our prompting  
 484 strategy with RAP can be found in appendix H.

485 We use the Python library `choix` (26) to find the maximum-likelihood estimate of agent ratings in the  
 486 Bradley–Terry model. This library is made available under the MIT License.

487 There is one extra game in the benchmark that can be found on the Github repository that was  
 488 not included in data collection: *Atari Boxing*. Data collection on this game turned out to be too  
 489 cumbersome, but as the only real-time game, it measures a factor not covered by the other games,  
 490 and thus is important for future benchmarking.

491 All games received two agents regardless of team size. In cases with multiple cooperative players on  
 492 one team, the agent is duplicated. The agent is not made explicitly aware that it is duplicated.

493 All code for running the benchmark on existing models and scaffolds, for creating implementing  
 494 new agents, and for reproducing results can be found in our Github <https://anonymous.4open.science/r/GameBench-5942/>

496 **H RAP prompting**

497 Reasoning-via-Planning describes a framework for using a probabilistic language model in a Monte  
 498 Carlo tree search. Exactly how the model is prompted depends on the implementation. Below, the  
 499 [rules] and [rules subtopics] come from Appendix K.

500 **Rules subtopics**

501 If you would like to learn more about the rules at any point, use rule(<topic>) where <topic> is one of [subtopics from rules above].

502 **!PREFIX**

503 You are now playing a game called [title]. The rules are as follows [rules summary from game above]. [Rules subtopics]. Your observation of the game is between <state> and </state>: <state>[game state]</state>

504 **!EXAMPLE**

505 You are playing a game called monty hall. The rules of the game are as follows: there are three doors, behind one of which there is a prize. Select the door with the prize. Your observation of the state is between <state> and </state>: <state>All three doors are closed.</state>

506 **Prompt for list of available actions**

!EXAMPLE  
User: To the best of your ability, predict the available actions in this position between <action> and </action>:  
507 Assistant: <actions> 0. Choose the left door 1. Choose the middle door 2. Choose the right door</actions>  
!PREFIX  
User: To the best of your ability, predict your available actions in this position between <actions> and </actions>:

508 **Prompt for selecting an action** The next-token probability from this prompt is used in the reward calculation.  
509

!PREFIX  
510 User: To the best of your ability, predict your available actions in this position between <actions> and </action>:  
Assistant: <actions>[actions from previous model prediction, enumerated]</actions>  
User: Choose an action by writing only the associated number.

511 **Prompt for self-evaluating an action** The next-token probability from this prompt is used in the reward calculation.  
512

513 User: Write your action below:  
Assistant: [action from previous model prediction]  
User: Is this a good action? yes/no.

514 **Prompt for guessing other players' actions**

!EXAMPLE  
515 User: To the best of your ability, predict what actions other players might take between <others> and </others>:  
Assistant: <others>My opponent is going to reveal one of the two doors I don't choose.</others>  
!PREFIX  
User: To the best of your ability, predict what actions other players might take between <others> and </others>:

516 **Prompt for guessing the game state**

```

!EXAMPLE
User: Write your action below:
Assistant: I will choose the left door
User: Write other player's actions below:
Assistant: My opponent will reveal the middle door
User: To the best of your ability, predict your new observation of the game based on your actions
and others' actions between <state> and </state>:
Assistant: <state>\nThe left and right doors are closed, and the middle is open. There is no prize
behind it.\n</state>
!PREFIX
User: Write your action below:
Assistant: [action from previous model prediction]
User: Write other players' actions below:
Assistant: [other players' actions from previous model predictions]
User: To the best of your ability, predict your new observation of the game based on your actions
and others' actions between <state> and </state>:

```

518 **Prompt for open-ended actions** The API we designed allows games to give “open-ended actions” to  
519 agents, in which they don’t select from a predefined list of options but instead provide a text response  
520 as the action. However, RAP doesn’t support this format, so we convert open-ended actions into  
521 ones with predefined options by prompting the model for a response to the open-ended action before  
522 feeding it into the Monte Carlo tree search algorithm.

```

!EXAMPLE
User: Write your action below:
Assistant: Ask my opponent a question.
User: This is an openended action. Write a description of what you're going to do.
Assistant: I will pretty-please ask them to tell me which door has the prize.
!PREFIX
User: Write your action below:
Assistant: [openended action from game]
User: This is an open-ended action. Write a description of what you're going to do.

```

524 **Prompt for assessing win probability** The next-token probability from this prompt is used in the  
525 reward calculation.

```

!PREFIX
User: Will you eventually win from this position? yes/no

```

## 527 I Additional game selection details

528 In order to evaluate a broad range of cognitive skills associated with strategic reasoning, we curated a  
529 diverse set of games featuring abstract strategy, non-deterministic outcomes, hidden information, lan-  
530 guage communication, social deduction and bluffing, and cooperation between players. A breakdown  
531 of which games had these features can be found in Table 3.

532 Using these categories, we then filtered for games unlikely to be significantly represented in LLMs’  
533 pretraining data, to evaluate the models’ out-of-distribution reasoning abilities. Two key criteria  
534 were (a) excluding games with dedicated online forums discussing improvement strategies, as well  
535 as (b) excluding games with published strategy guides. After finalizing the selection of games, we  
536 formalized their rulesets and mechanics into programmatic environments that the LLM agents could  
537 interact with.

538 Our final selection of games were *Air, Land, Sea* (ALS); *Arctic Scavengers* (ARC); *Are You the*  
539 *Traitor?* (AYT); *Codenames* (CN); *Hive* (HV); *Pit* (PT); *Santorini* (SN); *Two Rooms and a Boom*  
540 (TRB); and *Sea Battle* (SB). Descriptions of the games and their rules can be found in Appendices J  
541 and K respectively. For additional details about game implementation, see Appendix G.

Table 3: **Number of games per reasoning category** We identify a set of six orthogonal components of strategic reasoning and curate a set of games that sufficiently cover their spread.

Reasoning Category	Total	Games
Abstract Strategy	6	ALS, ARC, CN, HV, SN, SB
Non-Deterministic	3	ARC, TRB, SB
Hidden Information	3	ARC, AYT, TRB
Language Communication	4	AYT, CN, PT, TRB
Social Deduction	2	AYT, TRB
Cooperation	4	AYT, CN, SB, TRB

## 542 J Game descriptions

543 **Air, Land, and Sea** is a war strategy game where players are Supreme Commanders fighting to  
 544 control two of three areas (air, land, sea) by deploying limited Battle card forces each round. The first  
 545 commander to accumulate 12 points across multiple battles wins the war. <https://boardgamegeek.com/boardgame/247367/air-land-and-sea>

547 **Arctic Scavengers** is a resource-management game in which players are the leader of a small tribe  
 548 of survivors. Resources, tools, medicine, and mercenaries are all in scarce supply. Players are  
 549 pitted against each other in a fight for survival. The agent with the largest tribe at the end of the  
 550 game is declared the winner and receives 1 point. <https://www.riograndegames.com/games/arctic-scavengers-with-recon-expansion/>

552 **Are You the Traitor** is a social deduction game where players are secretly divided into Good and  
 553 Evil teams. The players then engage in an unstructured conversation trying to deduce the opposing  
 554 team’s critical roles. A player will yell “Stop!” while pointing at someone, and that round ends. If  
 555 they identify their target role correctly, their team earns Treasure cards. The team with the most  
 556 Treasure after multiple rounds wins. <https://www.looneylabs.com/games/are-you-traitor>

557 **Codenames** is a 2v2 cooperative game with one spymaster and one operative per team. All players  
 558 see a grid of words, and it is the spymasters’ job to create one-word clues that relate to multiple  
 559 predetermined words from the grid at once, and operatives must keep using these clues to guess  
 560 all of their team’s words. Agents are awarded more points for correctly guessing more words.  
 561 <https://boardgamegeek.com/boardgame/178900/codenames>

562 **Hive** is a strategy game occurring on a hexagonal grid. Each player has a team of bugs, each with  
 563 a unique skillset. Players try to coordinate their bugs in order to completely surround the enemy’s  
 564 queen bee. The winning agent is awarded 1 point. <https://www.gen42.com/games/hive>

565 **Pit** is an every-person-for-themselves trading simulation. Each player has a hand of cards, and each  
 566 card represents a certain commodity in the market. Players must trade semi-blindly trade cards to try  
 567 to obtain enough of any commodity to “corner the market.” Agents are awarded points based on the  
 568 commodity that they corner the market with. <https://www.gamenightgames.com/win1012.html>

569 **Santorini** is a strategy game in which two players take turns moving one of their two pawns on a five  
 570 by five grid and building blocks on the grid. The game ends when one of the players moves a pawn to  
 571 a square that has been built three blocks high or when one of the players cannot make a move. The  
 572 winning agent is awarded 1 point. <https://roxley.com/products/santorini>

573 **Two Rooms and a Boom** is a cooperative social-deduction game in which all players are split  
 574 into two teams and then mixed around between two rooms. No two players start knowing other  
 575 players’ teams or roles on the team, but it is the red team’s goal to end the game with the red-  
 576 team bomber and blue-team president in the same room, and it is the blue team’s goal for the  
 577 opposite. The winning agent is awarded 1 point for satisfying their team’s objective. <https://www.tuesdayknightgames.com/products/two-rooms-and-a-boom>



579 **Sea Battle** is 3v3 board game in which players' attempt to sink their opponents' ships and their  
580 movement and cannon-firing actions occur simultaneously. The winning agent is awarded 1 point  
581 if they eliminate all enemy ships before themselves becoming eliminated. [https://yppedia.  
582 puzzlepirates.com/Sea\\_battle](https://yppedia.puzzlepirates.com/Sea_battle)

## 583 **K Game rules**

584 The rules as follows are exactly as they were shown to the language models. Rules in bullet points  
585 were withheld until requested by a model taking a specific action "Explain(*rule heading*)".

586 **Arctic Scavengers** The game is played in 6 rounds, with each round consisting of a resource gathering  
587 phase and a skirmish phase. In the resource gathering phase, players draw cards from their deck  
588 and take actions to gather resources from the mercenary piles and the junkyard pile. In the skirmish  
589 phase, players compare the strength of their tribes and the winner of the skirmish gains a contested  
590 resource card. The game ends when all contested resource cards have been won, and the player with  
591 the largest tribe is the winner.

592 **Are you the traitor?** The Good team wants to destroy an Evil Magic Key while the Evil team wants  
593 to keep it. The key can be destroyed by giving it to the Good Wizard, but there is an Evil Wizard who  
594 looks exactly alike. Use social deduction to find out who is who, but also know that there is a traitor  
595 among the guards who have the key.

596 **Two Rooms and a Boom** Two teams, Blue and Red, have opposing goals. At the end of three  
597 rounds the Red team wants to have both the President and the Bomber in the same room, while Blue  
598 team wants them to be in opposite rooms. Each round will allow the Leader of each room to trade  
599 'hostages' in order to find out who the President and Bomber are and use that info to achieve their  
600 team's mission. Find out information by talking to other hostages in your room.

601 **Air Land and Sea** A strategic card game where two players compete over a series of battles to control  
602 different Theaters of war: Air, Land, and Sea. Each player is dealt 6 cards representing various  
603 military units and tactics. Players win a battle by controlling more Theaters than their opponent or  
604 convincing their opponent to withdraw. Victory Points (VPs) are earned by winning battles, and the  
605 first player to reach 12 VPs wins the game. Players must carefully manage their hand and strategically  
606 deploy cards to outmaneuver their opponent.

607 • **Battle Structure** During a Battle, the players take turns playing one card at a time, trying to  
608 control more Theaters than their opponent. You don't draw cards during a Battle, so be sure  
609 to plan carefully and make the most of the 6 cards you are dealt!

610 • **Theaters** Each of the three Theater boards creates a 'column' between the players: one for  
611 Air, one for Land, and one for Sea. These columns are called Theaters. Cards are always  
612 played into these three Theaters. If a card is in a particular Theater's column, we say that the  
613 card is 'in that Theater.' Theaters that are next to each other are called 'adjacent Theaters.' A  
614 player owns all of the cards on their side of the Theater boards. During your turn, you will  
615 play cards only on your side of the Theaters.

616 • **Battle Cards** Cards are played to advance your war effort and how they are played will  
617 ultimately determine who wins the war (the game). Strength: Each card has a Strength  
618 value. If the total Strength of all the cards on your side of the Theater is higher than the total  
619 Strength of all the cards on your opponent's side of that Theater, you 'control' that Theater.  
620 Tactical Abilities: Most cards have a Tactical Ability along with Strength, which takes effect  
621 as soon as the card is played 'face up' to a Theater. These abilities are either 'Instant' or  
622 'Ongoing.'

623 • **Type of Battle Cards** There are three types of cards: 'Air,' 'Land,' and 'Sea' cards, which  
624 relate to the three Theaters. Normally, you may only play a card 'face up' to its matching  
625 Theater: Air cards in the Air Theater, and so on.

- 626 • **Facedown Cards** Cards can also be played 'facedown' as a 'wild card' in any Theater.  
627 Facedown cards always have a Strength of 2. 'Facedown' cards do not have any Tactical  
628 Abilities. You may see your own facedown cards at any time, but you may not see your  
629 opponent's 'facedown' cards.
- 630 • **Covered Cards** When a card is played to a Theater that already contains cards, the newly  
631 played card is placed so that it overlaps the previously played card, while still showing the  
632 top portion of it. Any card overlapped by another is called a 'covered card.' Similarly, any  
633 card that is not overlapped by another card is referred to as 'uncovered.'
- 634 • **Resolving Battle** During a Battle, players take turns starting with the player who has the 1st  
635 Player me Commander card. On your turn, you must take only one of these three actions:  
636 Deploy, Improvise, Withdraw. Once you have finished your action, your opponent begins  
637 their turn. The players continue to alternate taking turns until one of them withdraws or both  
638 players have played all of their cards.
- 639 • **Possible actions:** Deploy: Play one card from your hand, 'face up.' When you play a card,  
640 you must follow these deployment restrictions: You can only play cards on your side of the  
641 Theater boards. The card must be the same type as the Theater you play it to. If you have  
642 other cards in that Theater already, you must place the new card so that it covers (partially  
643 overlaps) those cards. Improvise: Play one card from your hand, 'facedown', to any Theater.  
644 'Facedown' cards are treated as 'wild cards' and can be played to any Theater regardless of  
645 which type they are. Withdraw: If you think your chances of winning the current Battle are  
646 low, you may withdraw. If you do, your opponent wins the Battle and gains VPs depending  
647 on how many cards are left in your hand. See the me Commander cards for more specific  
648 information.
- 649 • **me Commander Cards** Supreme Commander Cards: The 1st Player Supreme Commander  
650 wins tied Theaters and gains the following number of VPs based on the number of cards left  
651 in their opponent's hand if their opponent withdraws: 5+ cards = 2 VPs, 3-4 cards = 3 VPs,  
652 2 cards = 4 VPs, 0-1 cards = 6 VPs. The 2nd Player me Commander loses tied Theaters and  
653 gains the following number of VPs based on the number of cards left in their opponent's  
654 hand if their opponent withdraws: 4+ cards = 2 VPs, 2-3 cards = 3 VPs, 1 card = 4 VPs, 0  
655 cards = 6 VPs.
- 656 • **Tactical Abilities** Most cards have Tactical Abilities described on the card. When you play  
657 a card face up from your hand, or if a facedown card is flipped over, its Tactical Ability  
658 takes effect immediately. There are two kinds of Tactical Abilities: 'Instant' and 'Ongoing',  
659 indicated on the card. You must carry out the effects of a Tactical Ability unless they contain  
660 the word 'may'. If a Tactical Ability is impossible to perform, that ability is ignored and has  
661 no effect.
- 662 • **Instant Abilities** Instant Abilities take effect immediately after the card is played or if the  
663 card is revealed by being flipped face up. Once the Instant Ability is resolved, it has no  
664 further effect (unless somehow that card is played or revealed again). Note: Because instant  
665 abilities take effect when flipped face up, it is possible for multiple instant abilities to take  
666 effect around the same time. In these situations, always resolve the instant abilities in the  
667 order they happened and fully resolve each ability before moving on to the next. Once an  
668 instant ability begins taking effect, it always resolves fully, even if it gets flipped facedown  
669 before completing.
- 670 • **Ongoing Abilities** These are always in effect as long as the card is face up. If a card with an  
671 Ongoing Ability is flipped 'facedown', the ability no longer has any effect (unless that card  
672 is revealed again). Example: The Escalation Tactical Ability increases the Strength of all  
673 of your facedown cards to 4 as long as the Escalation card remains 'face up'. If that card  
674 were flipped over by another Tactical Ability, your 'facedown' cards would go back to being  
675 Strength 2.
- 676 • **Tactical Ability Key Terms** Flip: Many Tactical Abilities allow you to flip a card. Flipping  
677 a card means either turning it 'face up' if it is 'facedown' or turning a 'facedown' card

678 so it is 'face up.' Unless the ability states otherwise, you may flip any card; yours or your  
679 opponent's. Uncovered/Covered: Many Tactical Abilities only affect uncovered or covered  
680 cards. If an ability does not specify uncovered or covered, such as Transport or Redeploy,  
681 assume the ability can affect any card. Play: Some Tactical Abilities instruct you to play  
682 a card, or only take effect in response to a card being played. The word 'play' describes  
683 any time a player takes a card from their hand and places it in a Theater. Non-Matching  
684 Theaters: Means that a card is not in the Theater of its type. The card does not suffer any  
685 penalty for being in the 'wrong' Theater. Destroy: Some Tactical Abilities instruct you  
686 to destroy a card. Destroyed cards are always placed facedown on the bottom of the deck.  
687 If a card is destroyed immediately after it is played, such as by Blockade, then that card  
688 does not get to use its Tactical Ability. Occupied: When counting the number of cards that  
689 occupy a Theater, always count both players' cards towards that total. Move: When a card  
690 is moved to a different Theater. It stays on the same side of the Theaters it was already on  
691 and remains owned by the same player. Moved cards are placed on top of any cards already  
692 in the Theater it was moved to. It covers those cards.

693 • **Ending Battles** There are two ways a Battle can end: If a player withdraws, their opponent  
694 wins the Battle. Or if both players have played all of the cards in their hand. At this point,  
695 the player who controls the most Theaters wins the Battle. In order to control a Theater,  
696 you must have a higher total Strength there than your opponent has in that Theater. If your  
697 Strengths are tied, the 1st Player wins the tie and controls that Theater. If there are no cards  
698 on either side of the Theater, the 1st player controls that Theater.

699 • **Scoring Victory Points** If neither player withdraws, the winner of the Battle scores 6 VPs.  
700 If one of the players withdraws, the other player scores VPs based on the number of cards  
701 left in the withdrawing player's hand (see the me Commander Cards for details). After  
702 scoring VPs, check if the victor has enough VPs to win the game (12 VPs). If they don't,  
703 fight another Battle.

704 • **Setting up Battles** All cards are collected and shuffled together to create a new deck. Deal  
705 each player a new hand of 6 cards. Next, the Theater cards are rotated clockwise so that the  
706 rightmost Theater is moved to the far left of the Theater lineup. Lastly, players swap me  
707 Commander cards. The player who was 1st in the last battle is now 2nd.

708 **Codenames** A strategic game of guessing and deduction where two teams, Red and Blue, compete to  
709 identify their team's words on a grid based on one-word clues given by their Spymasters. The game  
710 ends when all words of one team are guessed, or the assassin word is chosen.

711 • **Roles** Spymaster: Knows which words correspond to which team / the assassin. Gives  
712 one-word clues that relate to any number of their team's words on the board. Operative:  
713 Guesses words belonging to their team based on the Spymaster's clues. Aims to avoid words  
714 not belonging to their team and the assassin word.

715 • **Turn Structure** Spymaster's Turn: Give a clue to their operative and a number indicating  
716 how many words relate to that clue. Operative's Turn: Guess words, aiming to find all their  
717 team's words. After each guess, if the word is not their team's, the turn ends. If the word  
718 is their team's, they can guess again. If the word is the assassin word, the game ends and  
719 their team loses. An operative can make up to N+1 guesses, where N is the number of cards  
720 given by the Spymaster.

721 • **Winning Conditions** A team wins by correctly guessing out all their words. Game ends  
722 immediately if the assassin word is guessed and the team who guessed it loses.

723 • **Forbidden Actions** Spymasters cannot use part or any form of the words on the board in  
724 their clues. Spymasters cannot use words that sound like words on the board in their clues.  
725 Clues must be exactly one word and one number.

726 • **Scoring** Points are awarded based on the number of correct guesses by each team. If a team  
727 guesses the assassin word, they receive a score of 0.

728 • **Special Rules** If zero words are related to the clue, the Spymaster can give a clue of '0' and  
729 the Operative can guess an unlimited number of words.

730 **Hive** Hive is a bug-themed abstract strategy game. The object of Hive is to capture the opponent's  
731 queen bee by allowing it to become completely surrounded by pieces belonging to either player,  
732 while avoiding the capture of one's own queen. Tiles can be moved to other positions after being  
733 placed according to various rules, much like chess pieces.

734 • **Placing the Queen Bee** Players must place their Queen Bee by their fourth turn. Until then,  
735 they cannot move any placed pieces.

736 • **Queen Bee Movement** The Queen Bee can only move one space at a time around the hive.

737 • **Spider Movement** The Spider can move exactly three spaces.

738 • **Ant Movement** Able to move to any empty space around the hive as long as other movement  
739 rules are not violated.

740 • **Grasshopper Movement** The Grasshopper can jump over over adjacent pieces, landing on  
741 the first empty space.

742 • **One Hive Rule** The tiles must always be connected; you cannot move a piece if it would  
743 break the hive into separate groups.

744 • **Freedom to Move** A piece can only move if it can physically slide to its new position  
745 without disturbing other tiles.

746 • **Max Turns** The game ends after 250 turns.If no Queen Bee is surrounded by the end of the  
747 game, the game is a draw.

748 **Santorini** Win by moving one of your pawns to the third level of the board or forcing your opponent  
749 to be unable to finish their turn. The game is played on a five by five grid of squares, and each player  
750 controls two pawns. Play alternates between the players, starting with player 1. The pawn that a  
751 player plays with alternates during each of their turns: for example, player 1 plays pawn A on their  
752 first turn, pawn B on their next turn, then pawn A, and so on. Blocks can be placed on squares on the  
753 board up to four blocks high, creating four possible height levels.

754 The board begins with no blocks placed, so every square begins at level 0. Before the game starts,  
755 each of the players takes turns placing each of their pawns on the board. A square is occupied if a  
756 pawn is on it.

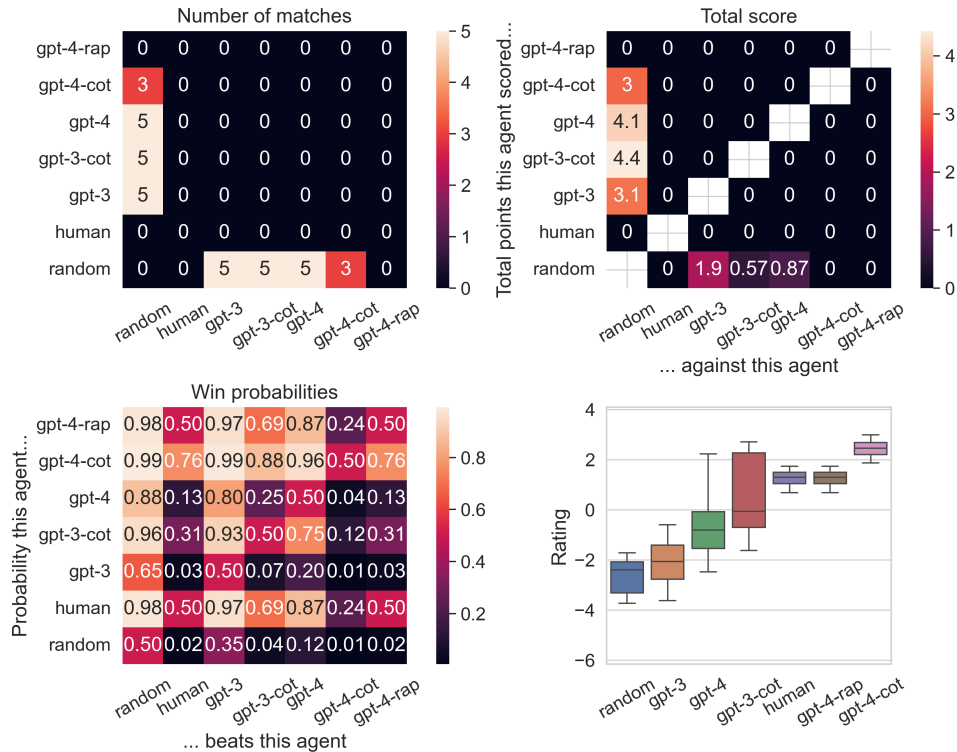
757 Each turn consists of two stages: the "move" stage and the "build" stage. During the move stage, the  
758 player moves their pawn by one square (horizontally, vertically, or diagonally). They cannot move  
759 their pawn onto a square that is occupied by another pawn, more than one level higher than the pawn,  
760 or at level 4. They can move a pawn any number of levels down, to the same level, or one level higher,  
761 but not more than one level higher and not to level 4.

762 During the build stage, the player must select an unoccupied square adjacent to the pawn they moved  
763 during the move stage and place a block on it. They can place a block onto an unoccupied square at  
764 any level less than 4. Once a square has been built to level 4, it is "complete", meaning pawns cannot  
765 move to it and blocks cannot be placed on it. The player instantly wins if they move their pawn onto  
766 a square at level 3 or if they force their opponent to not be able to finish their turn.

767 **Pit** Pit is a commodity trading game where players engage in trading to accumulate points and emerge  
768 as the winner. The game involves commodity cards representing various goods, with each card  
769 holding a specific point value. Players shout out their trade offers, attempting to negotiate deals with  
770 others to acquire valuable commodities. Additionally, Bull and Bear cards periodically influence the  
771 market conditions, either boosting or decreasing commodity values. The game continues with trading  
772 phases, market fluctuations, and scoring until a player or team reaches the agreed-upon point total,  
773 declaring them the victor in the spirited world of commodity trading.

774 **Sea Battle** Sink all of your opponent team's ships before they sink all of your team's ships.

## Air, Land, and Sea

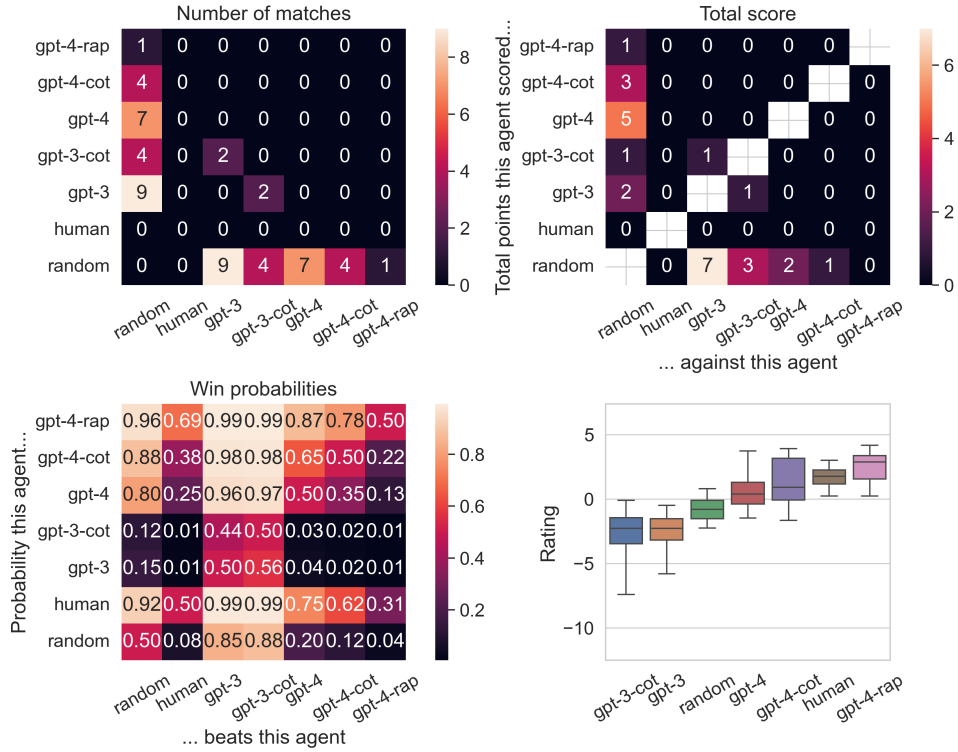


- 775
- 776
- 777
- 778
- 779
- 780
- 781
- 782
- 783
- 784
- 785
- **Damage** Players can be damaged in three ways: (1) by getting shot at by another player, (2) by sailing into a rock, (3) by colliding with another ship.
  - **Sinking** After a player has sustained enough damage, they sink and cannot play the rest of the round.
  - **Winning** A team wins if they have at least one live ship when all of their opponents have sunken.
  - **Board** The board is a 24x24 grid. Some squares are occupied by rocks and some are occupied by players' ships.
  - **Gameplay** Each turn, all players choose how they want to move and how they want to shoot. All players' choices are executed simultaneously.
  - **Teams** At the start of the game, there are three players on each team.

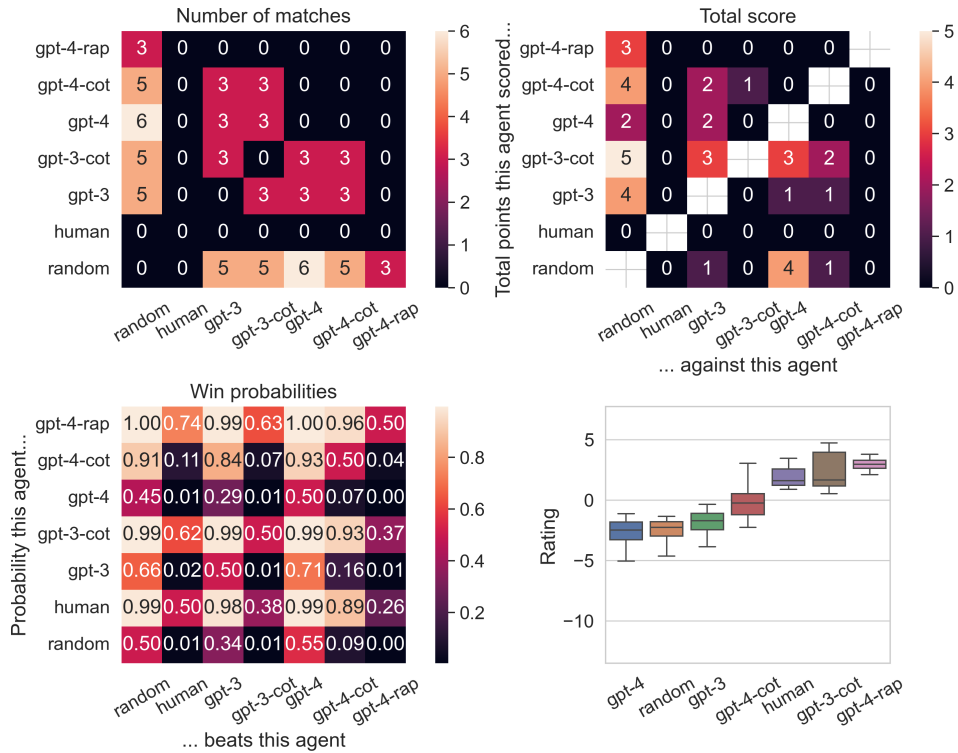
## 786 L Additional figures

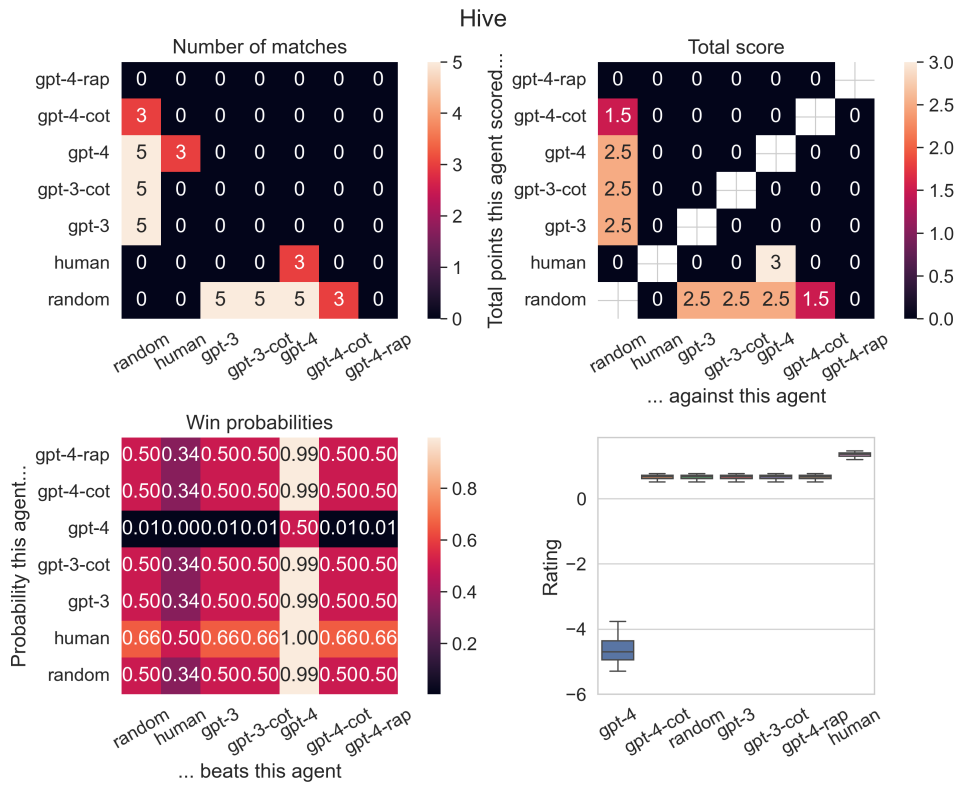
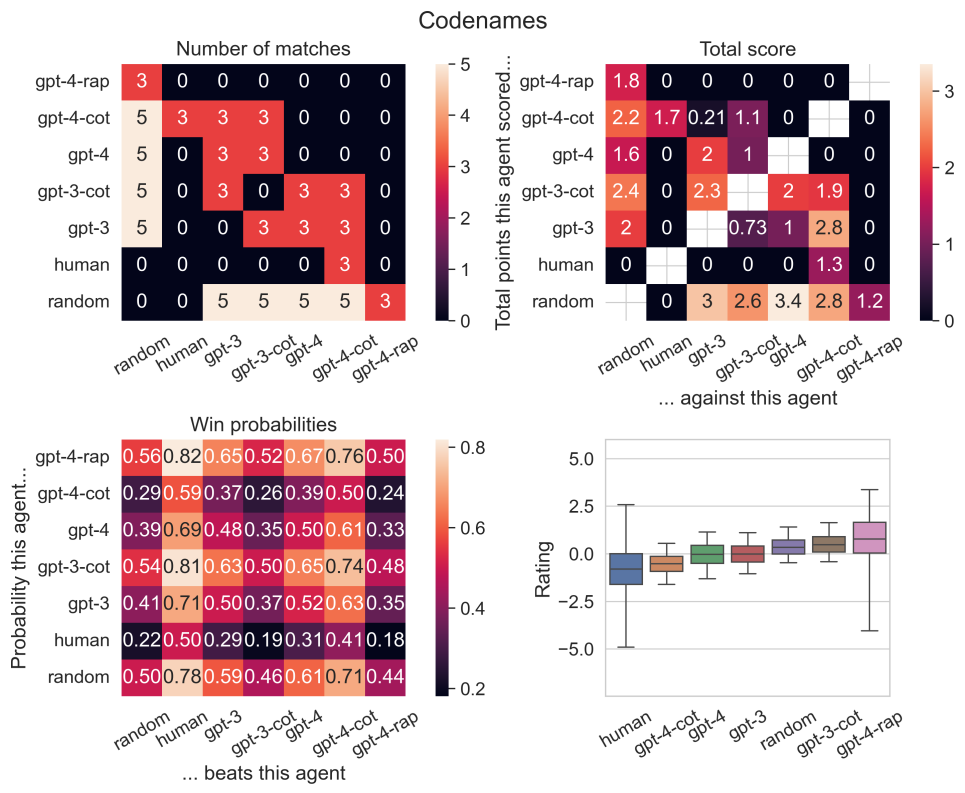
787 We present the match outcomes per game, including the number of matches, total score, win  
788 probabilities and rating per agent.

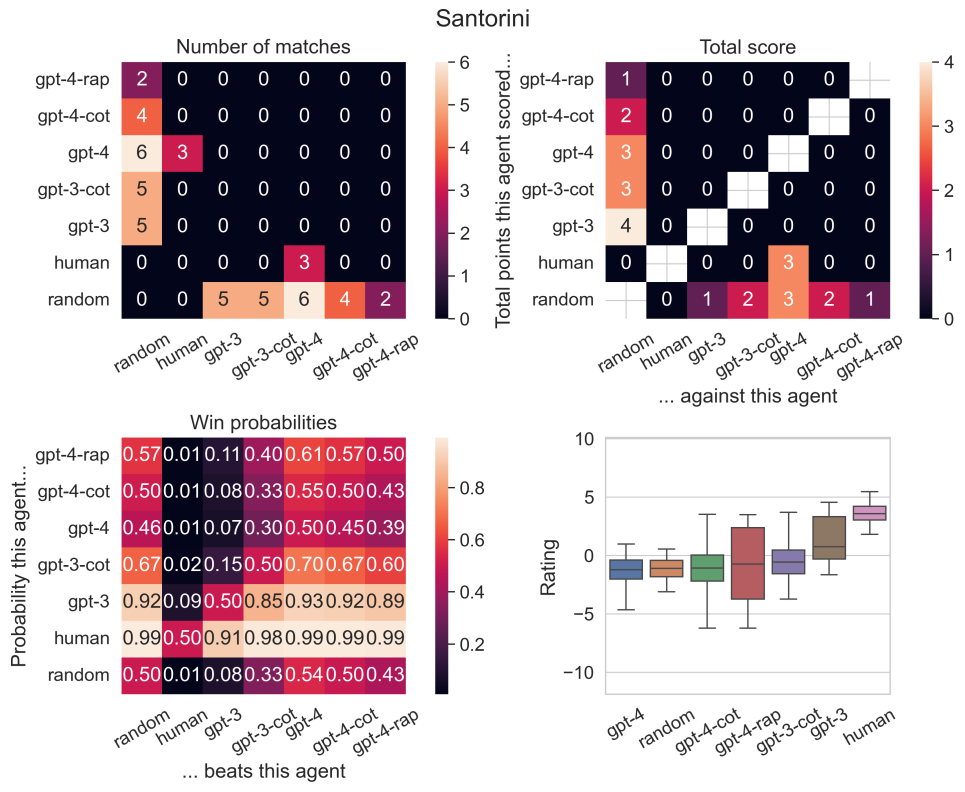
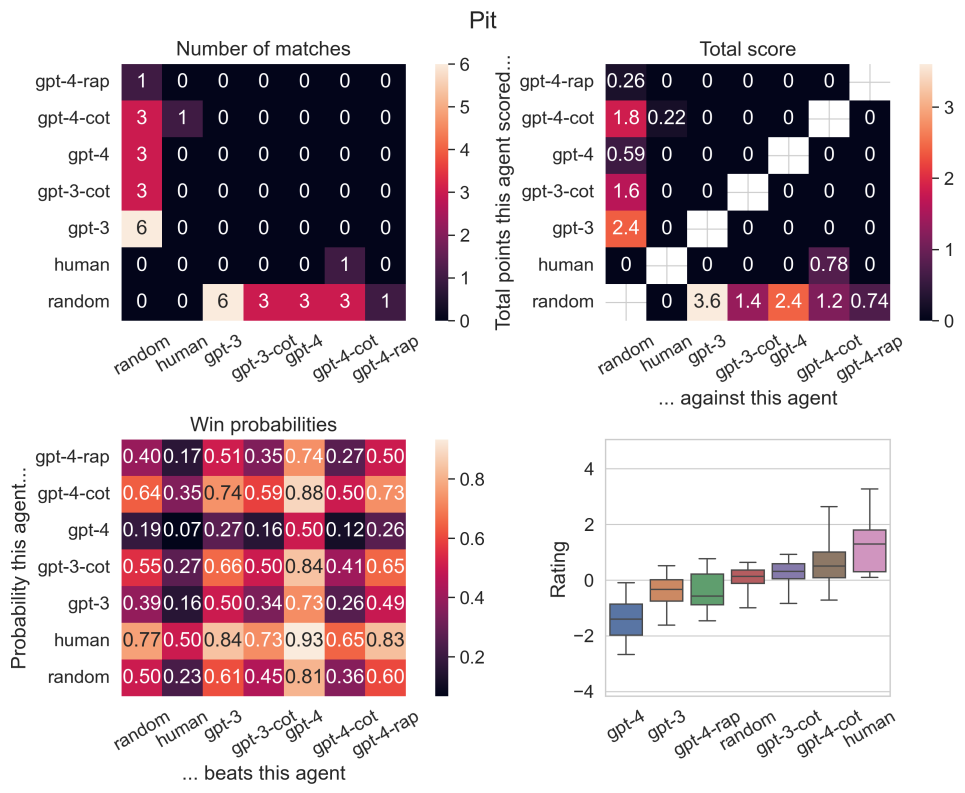
### Arctic Scavengers



### Are You the Traitor?

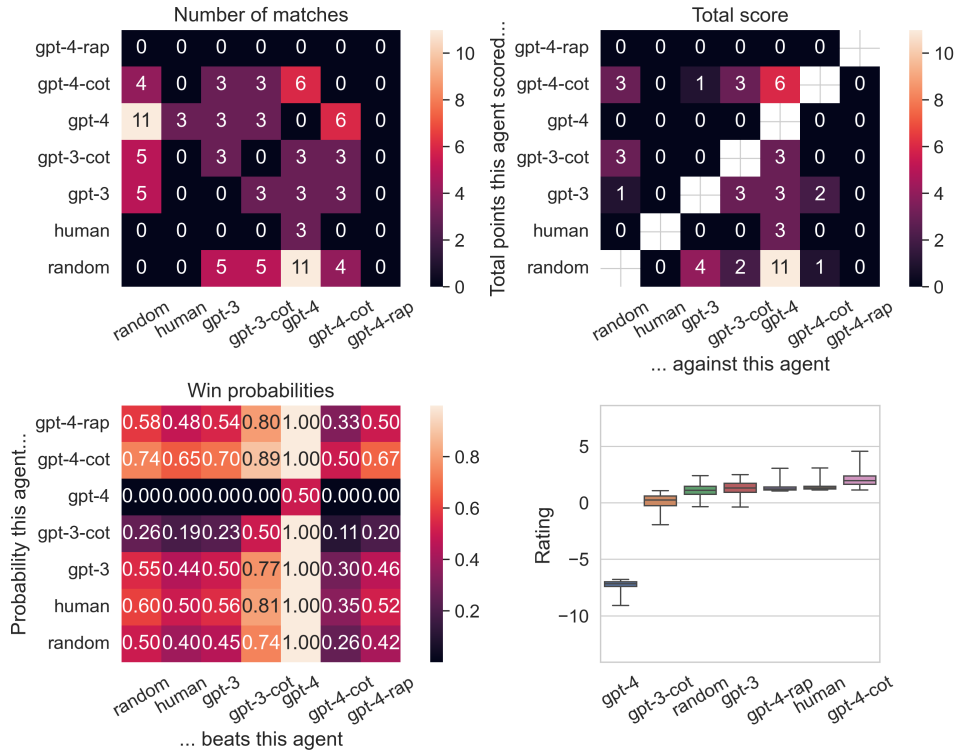








### Sea Battle



### Two Rooms and a Boom

