SPARKS OF COOPERATIVE REASONING: LLMs AS STRATEGIC HANABI AGENTS

Anonymous authorsPaper under double-blind review

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

033

035

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Cooperative reasoning under incomplete information is a significant challenge for both humans and multi-agent systems. The card game Hanabi ** embodies this challenge, demanding theory of mind reasoning and strategic communication. We present the largest evaluation to date of Large Language Models (LLMs) as Hanabi playing agents, assessing 17 state-of-the-art LLMs in 2 to 5-player cooperative multi-agent settings. We investigate why multi-agent coordination failures persist by systematically evaluating the impact of context engineering, from simple game state (Watson) tracking to scaffolding reasoning with explicit card deductions motivated by Bayesian inference (Sherlock) across a wide range of LLM capability (from 4B to 600B+ parameters). To our knowledge for the first time, we show 1) agents can maintain a working memory to track game state (Mycroft) instead of being explicitly provided engine deductions 2) a smooth interpolation of cross-play performance between different LLMs. In the Sherlock setting, the strongest reasoning models exceed 15 points out of 25 on average across all player counts, yet they still trail experienced human players and specialist Hanabi agents, both of which consistently score above 20. Lastly, we release the first public Hanabi datasets with move utilities and annotated game trajectories: 1) **HanabiLogs**: 1,520 full game logs for instruction tuning and 2) **HanabiRewards**: 560 games with dense move-level value annotations (rewards) for all candidate moves. Via instruction tuning on HanabiLogs, we show a 21% average score improvement with Owen3-4B-Instruct in the Sherlock setting, outperforming powerful closed-source LLMs like GPT-4o, Claude Sonnet 3.7 and Grok-3.

1 Introduction

Large Language Models (LLMs) have demonstrated significant success on tasks requiring complex individual ("single agent") reasoning, such as mathematics Lewkowycz et al. (2022), recently achieving gold medal performance at the 2025 International Mathematical Olympiad OpenAI (2025b); Luong & Lockhart (2025), and code generation Chen et al. (2021), with models now placing second at the AtCoder World Tour Finals OpenAI (2025a). However, a critical frontier lies in evaluating their ability to reason cooperatively. Recent benchmarks exploring interactive environments for LLMs often emphasize single-agent decision-making Hu et al. (2025a) or competitive dynamics Hu et al. (2025b). These settings do not adequately test the skills central to **cooperation**. Cooperative reasoning is essential for robust multi-agent systems and effective human-AI collaboration in real-world settings Mu et al. (2024), like coordinating autonomous vehicles in an intersection Liu et al. (2025) or collaborative robots on a factory floor. These settings involve interpreting ambiguous social cues from other, inferring hidden intentions from sparse signals, and coordinating decisions under uncertainty, and extend beyond single agent problem-solving skills.

To address this gap, we turn to *Hanabi*, a cooperative card game widely recognized for evaluating multi-agent reasoning and theory of mind Bard et al. (2020). In Hanabi, players are unable to see their own cards and must instead rely on limited communication and inference about other players' knowledge. Consequently, players must *continuously model their teammates' beliefs and intentions based solely on observed actions*, making Hanabi an ideal and challenging benchmark for cooperative strategy (for more on why Hanabi is an ideal benchmark, see Appendix A).

In this work, we evaluate the capability of state-of-the-art LLMs to cooperatively reason as multiagent Hanabi players. To establish baseline performance, we first provide agents with minimal context (MinCon), i.e. game state, legal moves, and simple instructions. To evaluate if agents can deduce information from prior teammate actions and the dynamic game state, we then equip each agent with deductive context (DeductCon), i.e. strategic advice and deductions about each teammate's hand based on previous clues (a form of game history). We alternatively refer to these as Watson (simple) and Sherlock (deductive) prompts, reflecting their relative reasoning capabilities.

We summarize our contributions as follows:

- 1. The **largest empirical evaluation** to date of **multi-agent cooperation of LLMs** through the lens of **Hanabi** (Section 3) in two to five player settings with 10 game seeds per setting.
- 2. An exhaustive investigation of the **factors that best equip LLMs with cooperative reasoning**, namely context engineering (Section 4), cross-play (Section 6.1), sampling methods (Section E.1), and specialized multi-agent scaffolding (Section E.2).
- 3. We introduce, to our knowledge for the first time, a multi-turn strategy that mimics human game state tracking that **requires agents to implicitly deduce information from the evolving game history** rather than rely on game engine-provided deductions (Section 5)
- 4. We release the **first public Hanabi datasets** with **move-level value estimates** and **annotated game trajectories**, HanabiLogs and HanabiRewards (Table 1). We show an average 21% score increase when training a lightweight LLM¹ on HanabiLogs (Section 6.2).

2 Related Work

 LLMs are increasingly evaluated in interactive settings that require planning, communication, and adaptive coordination, with recent work spanning cooperative games Wu et al. (2024), multi-agent environments Ma et al. (2024), and reasoning benchmarks Yang et al. (2024). The cooperative card game Hanabi has emerged as a particularly challenging testbed, widely regarded as a grand challenge for theory of mind reasoning and cooperation Bard et al. (2020). Early reinforcement learning (RL) approaches, including Bayesian Action Decoder (BAD), Simplified Action Decoder (SAD), and Off-Belief Learning (OBL) achieved scores of approximately 24/25 in a two-player setting with self-play, but performance degraded substantially for larger player counts and when paired with unfamiliar partners Hu et al. (2020; 2021).

Specialized RL policies for Hanabi Canaan et al. (2020) have recently been replaced with LLM agents, such as in LLM-Arena Chen et al. (2024) and SPIN-Bench Yao et al. (2024). However, LLM-Arena did not evaluate reasoning LLMs DeepSeek-AI et al. (2025), which show significant gains over instruction-tuned LLMs (Section 4). In contrast, SPIN-Bench includes recent reasoning LLMs but lacks a detailed study into the cooperative reasoning behind LLM decision-making for Hanabi as it focuses on wider evaluation coverage of different games and tasks. It also omits important experimental details such as the number of games or random seeds evaluated, making it difficult to replicate or assess the robustness of its findings. For example, SPIN-Bench shows a surprisingly low 6/25 two-player score for DeepSeek R1DeepSeek-AI et al. (2025) compared to 14.3/25 from our most basic setting, MinCon (Figure 4).

Targeted case studies have explored specific enhancement techniques for Hanabi. For example, Agashe et al. introduce a theory of mind reasoning step, followed by chain-of-thought prompting and answer verification to reduce fatal mistakes. Hybrid approaches such as Instructed RL Hu & Sadigh (2023) leverage LLMs to interpret human-written instructions and provide priors that guide smaller RL agents toward human-compatible conventions. Recently, Sudhakar et al. trained a text-based model (R3D2) to overcome the limitations of specialized Hanabi agents that struggle across different player counts, demonstrating that text-based Q-network learning can generalize to other player configurations. All of the above methods either embed a single LLM within a larger scaffold, evaluate only the 2-player setting, or rely on training a new model. In contrast, we evaluate 17 SoTA LLMs as Hanabi playing agents across 2 - 5 player settings with a progressive prompting schedule (Section 3).

¹Qwen3-4B-Instruct-2507

Table 1: A comparison of existing Hanabi datasets organized by their contributions towards number of games, player configurations, and annotations for move ratings and game trajectories.

Dataset	Games	Players		Move	Game
Dataset		Туре	Max Number	Ratings	Trajectories
HanabiData Eger & Others (2019)	1211	Human & Specialized Agent	2	Х	Х
AH2AC2 Dizdarevic et al. (2024)	3079	Human	3	×	X
HOAD Sarmasi et al. (2021)	4M	Specialized Agent	2	X	X
HanabiLogs (Ours)	1520	LLM Agent	5	Х	✓
HanabiRewards (Ours)	560	LLM Agent	5	✓	✓

We address three key limitations of existing work. Firstly, a *lack of transparency regarding* essential experimental details such as the number of games and seeds (Appendix G). This is especially important in Hanabi, where final scores are sensitive to initial conditions. A fair evaluation requires all agents to be assessed on the same set of seeds, and statistical significance requires multiple runs.

Secondly, existing evaluations are *not truly multi-turn*: they collapse cooperation into a single-prompt per turn that does not track game state in an agent's working memory. We therefore introduce a *multi-turn* setup (Section 5) that evaluates models' ability to cooperate by maintaining and updating their own state across turns, better reflecting real-world (human) gameplay.

Finally, to our knowledge, *no public dataset of move-level value estimates or large-scale, richly annotated game trajectories currently exists*, hampering reproducibility and advancement in RL-based post-training methods such as RL with verifiable rewards (RLVR) Lambert et al. (2024) and RL with AI Feedback (RLAIF). While several existing Hanabi corpora provide valuable resources, they remain incomplete for modern LLM research (see Table 1).

To address these limitations and ensure transparency and reproducibility, we provide complete details of our evaluation protocol, including the specific random seeds and number of games used for each configuration. We open-source game trajectories via **HanabiLogs**, which includes approximately 1,520 complete games covering 2 - 5 player counts; and **HanabiRewards**, which also contains dense move ratings for 560 games from reasoning LLMs. We hope that these contributions enable reproducible and fair benchmarking and provide a resource for post-training for cooperative reasoning.

3 EXPERIMENT SETUP

We utilize the Hanabi Learning Environment (HLE) Google DeepMind (2019) for our game setup. For each player (in our case, agent), HLE provides their explicit knowledge, i.e. what each player knows about their own cards; we provide this information in both Watson and Sherlock setups) and a list of possible colors and ranks for each card (provided only in the Sherlock setup), updated according to clues received. For instance, if a player holds a yellow 5 and receives a red clue, the possibility list for that card will exclude red. We visualize this explicit deductive context in Figure 1. For Sherlock, we also provide general Hanabi strategies, as well as step-by-step reasoning workflow inspired by Bayesian inference (See Section 3.2 and Appendix I.2) For details of our LLM evaluation suite, see Appendix B.

We evaluate agents across two, three, four, and five-player team settings. To ensure robust evaluation, each agent plays 10 games per setting using different random seeds, totaling 40 games per agent. All games are played with each player using the same LLM as a Hanabi playing agent, e.g. four GPT-4.1 agents playing as a four-player team. If a team loses all three life tokens, we record their score at the moment of failure, as is standard in prior benchmarks Yao et al. (2024); Chen et al. (2024).

3.1 Watson Setting

To allow agents to define their own gameplay and test their knowledge of Hanabi, we first provide agents with Minimal Context (MinCon / Watson). Each agent receives essential state variables: turn number, player number, available information and life tokens, and discard pile contents. The input also included visible cards in other players' hands and their inferred knowledge about their own

163

164

167

171

177 178 179

181 182

183

184

185

187

188

189

190

191

192

193 194

195 196

197

199

200 201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

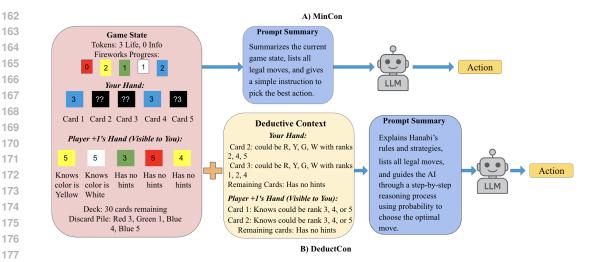


Figure 1: A comparison of the MinCon (Watson) and DeductCon (Sherlock) settings with an example 2-player Hanabi game state.

hands to assist clue selection (Figure 1 below Player +1's Yellow card 5 "Knows color is Yellow"). We found that omitting this perspective leads to agents giving redundant clues, as LLMs cannot infer what other players already know without a multi-turn trajectory. Agents are tasked with choosing the best move from a provided list of legal candidates, and also gave a rating (between -1 and 1) for each candidate, which we use to create the HanabiRewards dataset. All agent interactions, including reasoning traces from Qwen-3-225B-A22B, Qwen-3-32B, and Deepseek R1, are logged to compile our high-quality instruction tuning dataset, HanabiLogs. Once the deck is exhausted, we append "this is the final round and player+n is the last player" to the prompt. This ensures that agents are aware of the game's final round and can identify the last player to act, discouraging them from giving clues to players who would not have a turn and encouraging the last player to take risks rather than discarding or giving clues. We show an example of the Watson with o4-mini in Appendix I.1.

3.2 Sherlock Setting

We now focus on equipping agents with strategic reasoning by adding Deductive Context (DeductCon / Sherlock) to our agents. In our Sherlock setting, we use the Hanabi Learning Environment (HLE) Google DeepMind (2019) to provide explicit deductive feedback to the agent context Yao et al. (2024). We later discuss a variant of Sherlock where the agent must implicitly track its own deductive context over time (Section 5).

For example, as shown in Figure 1, the Deductive Context (the yellow box) specifies that "your card 2 could be Yellow, Red, Green, or White and Rank 2, 4 or 5", removing impossibilities based on prior clues (though discards are not considered in this deduction; agents must infer those independently). This approach provides agents with a snapshot of the game's trajectory from the game engine (HLE). To examine the effects of context engineering on Hanabi scores, we construct a systematic ablation study with the 5-player Seed 3 game using Grok-3-mini due to its favorable cost-performance trade-off, running each setup 10 times.

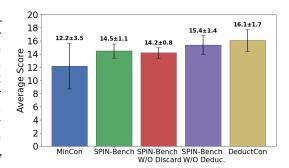


Figure 2: Average score with different prompt strategies for 10 runs of a 5-player game with Grok-3-mini. Error bars are standard deviation.

First, we compare our simple MinCon setting to SPIN-Bench, and observe a clear degradation of score from $14.5 \rightarrow 12.2^2$. Next, we evaluate the effect of providing card deductions to the agent by removing this additional information from SPIN-Bench. Specifically, we omit the "could be" possibilities for all players' hands (SPIN-Bench W/O Deduc). Surprisingly, agent performance slightly improved without these deductions ($14.5 \rightarrow 15.4$). This suggests that the agents did not effectively leverage deduction or discard-pile information to calculate probabilities. To further test this, we remove the discard pile from the prompt as well; performance slightly degrades (-0.3), but remains better than the MinCon setting (+2.0), indicating that the richer context or "prefill" the agent receives from SPIN-Bench over MinCon is generally beneficial.

Sherlock: Let's deduce step-by-step. To encourage the agent to actively use the additional deductive information provided, motivated by Bayesian inference, we ask the agent to calculate the probabilities for each card in its chain-of-thought before choosing its next action. We also include the starting card distribution and a final round flag similar to the Watson setting. As shown in Figure 2, Sherlock improves on the runner-up strategy (our deduction-less variant of SPIN-Bench) from $15.4 \rightarrow 16.1$. We provide all prompt variants in Appendix I.2.

4 BENCHMARK RESULTS

In this section, we benchmark the performance of Hanabi agents in our Watson and Sherlock settings and how performance varies across player counts. As shown in Figure 3, reasoning models, such as o3, o4-mini, Grok-3-mini, DeepSeek R1, Qwen-3-235B-A22B, Gemini 2.5 Pro/Flash, generally achieved higher scores (>13/25) than non-reasoning models (<10/25), even when game history information via deductions is not provided, i.e. the Watson setting. We find that reasoning models consistently benefit from deductive context provided by the Sherlock setting, with the exception of o4-mini in 4 and 5-player settings (see Figure 4). In contrast, adding Hanabi strategies and encouraging probabilistic reasoning (Sherlock) reduces performance in all non-reasoning models except Mistral Medium 3. We also find that deductive context (Sherlock) does not benefit all agents equally; while Gemini 2.5 Flash/Pro and Grok-3-mini improve substantially (+2.7 on average), o4-mini improves only slightly (+0.6 on average).

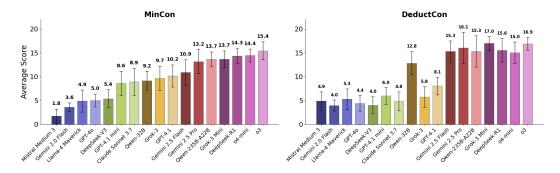


Figure 3: Score of 17 state-of-the-art LLM Hanabi agents averaged over two to five-player settings. We show scores for each specific player count in both settings in Figure 9 (Appendix D). Error bars denote standard deviation.

We show in Figure 4 that as player counts increase, Hanabi scores tend to drop for the best-performing reasoning LLM agents. DeepSeek-R1 (MinCon) and Gemini 2.5 Pro (DeductCon) are slight exceptions. We highlight that this performance drop is less severe than what has been reported by Sudhakar et al. for AI agents specifically trained for Hanabi (roughly $20+\rightarrow 15$ from 2-player to 5-player cross-play). This suggests that non-specialized LLMs acting as Hanabi agents may possess more robust and generalizable cooperative reasoning capabilities across different player counts compared to specialized agents.

²The high standard deviation for MinCon is due to a single early loss (score = 3/25). If we ignore this outlier, the mean score is 13.2, which is still 1 - 3 points less than all other strategies.

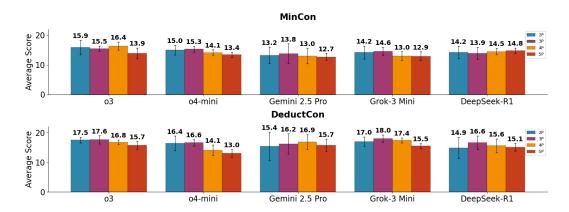


Figure 4: Average score of top-performing reasoning LLM based Hanabi agents when varying player count from 2 to 5. Error bars denote standard deviation.

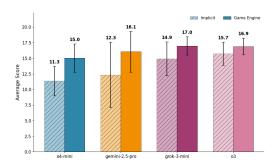
Excellent and Elementary: Watson vs. Sherlock. In the Watson (MinCon) setting, o3 outperformed all other agents for 2-4 players (Figure 4), but its scores dropped significantly in the 5player game, second to DeepSeek R1 (-0.9). In the Sherlock (DeductCon) setting, Grok-3-mini achieved the highest score for 3 (18.0) and 4 players (17.4), and only lagged behind o3 for 2 players (-0.5) and o3 and Gemini 2.5 Pro for 5 players (-0.2), showing consistently strong performance across player counts. Interestingly, we observe emergent strategies unique to each agent, even though they are provided the same context in each strategy: o4-mini discarded cards more frequently with Sherlock, whereas with the Watson prompt, it discarded only when out of information tokens. Gemini 2.5 Pro adopted an aggressive strategy until losing two life tokens, then shifted to conservative play. This sometimes led to the agent losing its last life token before the deck was exhausted. In contrast, Grok-3-mini consistently avoided losing life tokens, resulting in a low variance of scores compared to Gemini 2.5 Pro (Figure 4). Although the best reasoning models achieved average scores around 15-18 points out of 25, clearly surpassing earlier generations of LLMs, their performance remains below both state-of-the-art self-play search agents (>23 from Lerer et al.) and the recently introduced generalist Hanabi agent R3D2 (\geq 20 in 2, 3, and 4-player self-play; \approx 18 in 5-player setting from Sudhakar et al.). The agents' scores are also lower than those of experienced human Hanabi players (\sim 18-23), especially with few players (see Appendix F).

When changing context from Watson to Sherlock (Figure 3), among non-reasoning models, the GPT-4.1 family was relatively robust (-2.4 on average) compared to other agents, such as grok-3 (-3.9 on average) and Claude Sonnet 3.7 (-4.1 on average). For reasoning models, Gemini 2.5 showed comparable improvements with Sherlock(Flash: +4.4, Pro: +2.9). This provides some evidence for agents within a model family being similarly impacted by deductive reasoning enabled by providing richer contextual information (e.g. GPT 4.1, Gemini 2.5). We discuss more detailed turn analysis and agent behaviors in Appendix C.

Limitations of Sherlock. The primary limitation of the Sherlock setting is that we provide game history as explicit deductions from the Hanabi Learning Environment game engine (see Appendix I.2) rather than the agent implicitly deducing this information through its own interactions as the game progresses turn-by-turn. We attempted this multi-turn evaluation with a few agents, such as o4-mini and Grok-3 Mini, but were unable to run games longer than 30 turns due to LLM context window limits. We discuss a potential solution to this problem and introduce the multi-turn evaluation in Section 5.

5 Mycroft: Implicit Deductions from Multi-turn Play

Instead of providing the agent with programmatic deductions from a game engine (Sherlock), we hope to encourage the agent to implicitly deduce information from play thus far for their own future turns, similar to how a human would play the game. To this end, we provide agents information about their own cards only when a card has been directly clued (e.g., if a card is yellow and the agent receives a yellow clue, they know that card is yellow). We do not provide agents with other players'



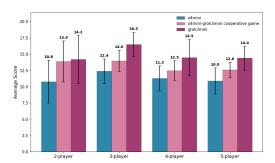


Figure 5: Average Hanabi scores for the best reasoning LLM agents with Implicit deductions (Mycroft) vs Game Engine deductions (Sherlock) averaged across 2-5 Player settings. Error bars denote standard deviation.

Figure 6: Average Hanabi score across 2–5 players for three team compositions: (left) all o4mini, (middle) one Grok-3-mini agent and the remaining o4-mini agents, and (right) all Grok-3-mini agents. Error bars denote standard deviation.

perspectives (e.g., "Player+1's Hand, Card1: Knows color is yellow," as shown in Figure 1) or with any deductive context about cards in any player's hand. Instead, agents are expected to infer such information themselves by reasoning over game history and to explicitly record their deductions, which is then made available to them on their next turn. Specifically, on any given turn, each agent's context includes the current game state and the agent's action in previous turn, serving as a working memory to track and update information across turns. To help the agent accurately update states, we instruct the agent on how the Hanabi Learning Environment (HLE) handles card positions after plays or discards (exact prompt in Appendix N). We term this setting Mycroft. In addition to cooperative reasoning, this setting also evaluates the agent's ability to deduce information by tracking its own behavior via multi-turn interaction over the game history, moving the needle closer to human strategy.

We evaluate the best performing reasoning LLMs from Sherlock setup which use engine-provided deductions, i.e. o3, o4-mini, Grok-3-mini and Gemini 2.5 Pro with the implicit deduction from multi-turn play (Mycroft). As shown in Figure 5, when Hanabi scores are averaged across player counts, o4-mini and Gemini 2.5 Pro consistently struggle to implicitly track the evolving game state based on the prior turn information, with a performance decline of \sim 3.7. Grok-3-mini shows a middling drop of \sim 2.1, while o3 shows the best multi-turn state tracking capability by dropping by only \sim 1.2. We provide detailed scores for each player setting (2 - 5) in Appendix N.2.

6 Ablations

6.1 Cross-Play

Thus far, we have only evaluated an LLM agent's ability to cooperate in self-play settings, i.e. when all players in a team are the same LLM (e.g. DeepSeek-R1). We now switch to a more realistic cross-play setting, where agents need to cooperate with teammates who are very different from them (i.e. other LLMs). We consider two LLMs with a wide performance gap in the Mycroftsetting, i.e. Grok 3 Mini (14.9) and o4-mini (11.3) to examine their ability to cooperate. Across 2-5 player settings, we always have exactly one Grok 3 Mini player, and the rest (1-4 players) are o4-mini. We choose this setting to examine whether adding a "stronger" player to weaker players makes the overall multi-agent system better.

As shown in Figure 6, the cross-play setting always performs better than the o4-mini self-play setting (Mycroft), and worse than the Grok 3 mini self-play setting (Mycroft). In other words, this provides preliminary evidence that cross-play performance interpolates between self-play performance of a weak and strong agent in the multi-turn setting.

6.2 INSTRUCTION TUNING ON HANABILOGS

Finally, to validate the effectiveness of our new datasets, we instruction tune Qwen-3-4B-Instruct-2507 on HanabiLogs (Section 2). We choose this LLM for its size and its strong instruction following and consistent output formatting capabilities, which is important when evaluating Hanabi. For a fair comparison against agents evaluated with Sherlock, we preserve the exact Sherlock prompt format. As our goal was to examine dataset capability for instruction tuning cooperative reasoning, we avoid "thinking" variants (also known as reasoning LLMs) to minimize conflating our dataset quality with gains provided by learned reasoning traces DeepSeek-AI et al. (2025).

Concretely, we instruction tune Qwen-3-4B-Instruct-2507 for 3 epochs on a subset of HanabiRewards containing only trajectories from o3 and Grok 3 Mini under the Sherlock setting, targeting imitation of their strong cooperative play (Figure 4).

As shown in Figure 7, the instruction tuned model improves by $\sim\!21\%$ and closes the gap with strong closed-source systems like GPT-40 and Grok-3. For example, post instruction tuning, Qwen3-4B reliably gives Rank-1 hints early and plays those cards, a behavior that was rare in the base model. We provide exemplar before/after comparisons via game transcripts in Appendix H.3.

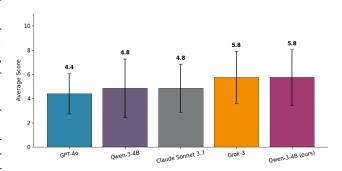


Figure 7: Average scores of Qwen-3-4B-Instruct-2507 before and after instruction tuning on **HanabiLogs** vs Grok, Claude Sonnet 3.7 and GPT-40. **Note**: We evaluated the Qwen-3-4B-Instruct-2507 models on different seeds to avoid memorization effects. Error bars denote standard deviation.

To further improve and match the performance of reasoning models, we suspect that agents may require a large increase in training compute (both data and GPU hours). We did not run reinforcement learning with verifiable rewards (RLVR) due to compute constraints; unlike math problems Luong & Lockhart (2025), learning a Hanabi policy with LLMs requires (implicitly) assigning probabilities over all candidate moves, with many rollouts and long horizons to discover good policies. We provide additional details, training hyperparameters, and discussion in Appendix H.

7 Future Work

Our high-level goal is to evaluate and improve the cooperative capabilities of LLMs in multi-agent settings, which we do in this work through the lens of Hanabi. A natural extension of this work is to evaluate if and how cooperative capabilities transfers across tasks, such as different games. For example, if we train an LLM with **HanabiRewards** (**Ours**), how well can the LLM transfer its cooperative ability to playing Overcooked Sun et al. (2025)?

We empirically show (Section 5) that even the best state-of-the-art reasoning LLMs still fall short at implicitly building their own deductions about game state from their own prior moves, which human players excel at. A natural extension of our multi-turn setup for implicit deduction is as an environment to train state tracking models with reinforcement learning, where the Hanabi Learning Environment (HLE) provides verifiable state tracking rewards. Since our dataset contains logs (HanabiLogs (Ours)) and dense move-level annotations (HanabiRewards (Ours)), we can use them together as one of the rollouts during RL training. Another valuable direction is to investigate how specialized training on games with verifiable rewards like Hanabi affects LLM generalization to other verifiable domains, such as mathematics and coding.

Lastly, our current cross-play setup (Section 6.1) only compares two LLMs in a single setting, i.e. progressively adding the stronger player to a single weaker player; there is significant scope for a more systematic and in-depth study of agentic cross-play for cooperative reasoning. This setting offers verifiable insight into real-world deployment scenarios, where multiple specialized agents that do specific tasks must all cooperate towards a higher goal. In our experiments, we observed that even when given identical instructions, different agents' strategies can diverge significantly (see

Section 4 and Appendix C). Recent works such as Dizdarevic et al. (2024) have made initial strides into Human-AI collaboration; we believe this direction is essential in developing more robust and adaptive cooperative AI systems.

8 CONCLUSION

In this work, we show via an exhaustive empirical evaluation of 17 state-of-the-art LLMs, including recent reasoning models, that while LLM agents show sparks of robust cooperative reasoning, they are **not yet fully generalist Hanabi agents**. The best performing reasoning LLMs (e.g. o3, Grok-3 Mini, Gemini 2.5 Pro) are limited in their ability to consistently infer teammate intentions and still fall short of both specialized Hanabi agents and strong human players (See Appendix F).

We propose two settings for cooperative reasoning (Figure 1, Section 3) one where we provide simplistic, minimal context to the agent (MinCon/Watson), and one where we provide Hanabi strategies and deductions from Hanabi game engine about player hands and enforce step-by-step probabilistic reasoning (DeductCon/Sherlock).

We empirically demonstrate that agents can generalize across different player counts (Section 4 and Appendix E.2, Figures 4, 11 and 12) and score reasonably well (>13/25) even when the games historical context is not explicitly provided by game engines (Watson), indicating that agents are not simply memorizing solutions for specific scenarios. When switching out explicit engine deductions (Sherlock) for encouraging the model to implicitly track state from its own previous turns (a novel task for Hanabi, which we call Mycroft), we empirically demonstrate that even state-of-the-art reasoning models like o3 and Grok-3 fail to accurately track game state, with an average performance decline of 2.7 (Section 5). We also show that using specialized-role agents is not a universal solution: in some scenarios, a well-steered simple agent (Sherlock) can perform equally well when provided with detailed context (Appendix E.1), and in some cases, prefilling the context of a mixture of specialized agents with diverse, relevant information helps (Appendix E.2).

When evaluating the capability of different LLMs to cooperate (cross-play), we observe sparks of cooperative reasoning: successively adding stronger players improves team performance; in 5 player settings. With only 1/5 of the turns played by the stronger agent, scores increase by ≈ 1.7 (Section 6.1). Our observed improvements from context engineering suggest that LLMs have untapped cooperative reasoning potential that could be further developed through improved training methods. To this end, we create the first public Hanabi datasets that have move-level value estimates and annotated game trajectories, HanabiLogs and HanabiRewards (Section 2). We empirically demonstrate the potential of our new dataset by instruction tuning a lightweight LLM, Qwen3-4B-Instruct on the o3 and Grok 3 Mini subset of HanabiLogs, which improves its score by 21% and surpasses strong closed-source LLMs like Grok-3, GPT-40, and Claude Sonnet 3.7 (Section 6.2).

9 REPRODUCIBILITY STATEMENT

We discuss the primary limitations of benchmarking cooperative multi-agent systems for Hanabi (Section 2), and highlight a lack of reproducibility due to missing experiment setup details in prior work. To this end, we provide the prompts for our agents, and crucial details such as the number of games and specific seeds for each setting in appendix I.1 I.2 G L M N. To further research in evaluating the cooperative reasoning capabilities of multi-agent systems via Hanabi, we commit to fully open sourcing our two new datasets, HanabiLogs and HanabiRewards (Section 2), after publication. We hope these datasets will prove valuable to the community for post-training cooperative reasoning of LLMs. Lastly, we commit to fully open sourcing **all** of our code and models trained on the HanabiLogs dataset after publication.

REFERENCES

- Saaket Agashe, Yue Fan, Anthony Reyna, and Xin Eric Wang. Evaluating multi-agent coordination abilities in large language models. *Findings of the Association for Computational Linguistics:* NAACL, 2025. URL https://arxiv.org/abs/2310.03903.
- Anthropic. Claude 3.7 sonnet system card. Technical report, Anthropic, 2025. URL https://www.anthropic.com/claude-3-7-sonnet-system-card.
- Nolan Bard, Jakob N. Foerster, Sarath Chandar, Neil Burch, Marc Lanctot, H. Francis Song, Emilio Parisotto, Vincent Dumoulin, Subhodeep Moitra, Edward Hughes, Iain Dunning, Shibl Mourad, Hugo Larochelle, Marc G. Bellemare, and Michael Bowling. The hanabi challenge: A new frontier for AI research. *Artificial Intelligence*, 280:103216, 2020. doi: 10.1016/j.artint.2020.103216.
- Antoine Bauza. Hanabi. Board game published by R&R Games, 2010. Spiel des Jahres winner 2013, https://en.wikipedia.org/wiki/Hanabi_(card_game).
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with gpt-4. https://arxiv.org/abs/2303.12712, 2023.
- Rodrigo Canaan, Xianbo Gao, Youjin Chung, Julian Togelius, Andy Nealen, and Stefan Menzel. Behavioral evaluation of hanabi rainbow dqn agents and rule-based agents. *Sixteenth AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, 16(1):31–37, 2020.
- 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. https://arxiv.org/abs/2402.16499, 2024.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. https://arxiv.org/abs/2107.03374, 2021.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, and et al. Lintz, Nathan. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025. URL https://arxiv.org/abs/2507.06261.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J.L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R.J. Chen, R.L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S.S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W.L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun

541

543

544

546

547

548

549

550

551

552

553

554

558

559

561

562

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580 581

582

583 584

585

586

587

588 589

590

592

Gao, Wenqin Yu, Wentao Zhang, X.Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y.K. Li, Y.Q. Wang, Y.X. Wei, Y.X. Zhu, Yang Zhang, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z.F. Wu, Z.Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2024. URL https://arxiv.org/abs/2412.19437v1.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z.F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J.L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R.J. Chen, R.L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S.S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wengin Yu, Wentao Zhang, W.L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X.Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y.K. Li, Y.Q. Wang, Y.X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y.X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z.Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. arXiv preprint arXiv:2501.12948, January 2025. URL https://arxiv.org/abs/2501.12948.

- Tin Dizdarevic, Tobias Gessler, Ani Calinescu, Jonathan Cook, Matteo Gallici, Andrei Lupu, and Jakob Nicolaus Foerster. Ah2ac2: A dataset for ad-hoc human—ai coordination in hanabi. https://openreview.net/pdf?id=Kioojohsuy, 2024.
- Markus Eger and Others. Hanabidata: A dataset for cooperative card game research. https://github.com/yawgmoth/HanabiData/tree/paper, 2019.
- Google DeepMind. The Hanabi Learning Environment. https://github.com/deepmind/hanabi-learning-environment, 2019.
- Google DeepMind. Introducing gemini 2.0: our new ai model for the agentic era. https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/, 2024.
- Google DeepMind. Start building with gemini 2.5 flash. https://developers.googleblog.com/en/start-building-with-gemini-25-flash/, 2025.

- Hengyuan Hu and Dorsa Sadigh. Language instructed reinforcement learning for human-ai coordination. https://arxiv.org/abs/2304.07297, 2023.
 - Hengyuan Hu, Adam Lerer, Alex Peysakhovich, and Jakob N. Foerster. "other-play" for zero-shot coordination. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020. URL https://arxiv.org/abs/2003.02979.
 - Hengyuan Hu, Adam Lerer, Brandon Cui, David Wu, Luis Pineda, Noam Brown, and Jakob Foerster. Off-belief learning. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2021. URL https://arxiv.org/abs/2103.04000.
 - Lanxiang Hu, Mingjia Huo, Yuxuan Zhang, Haoyang Yu, Eric P. Xing, Ion Stoica, Tajana Rosing, Haojian Jin, and Hao Zhang. Imgame-bench: How good are llms at playing games? *arXiv preprint arXiv:2505.15146*, 2025a. doi: 10.48550/arXiv.2505.15146. URL https://arxiv.org/abs/2505.15146.
 - Lanxiang Hu, Qiyu Li, Anze Xie, Nan Jiang, Ion Stoica, Haojian Jin, and Hao Zhang. Gamearena: Evaluating Ilm reasoning through live computer games. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025b. URL https://arxiv.org/abs/2412.06394.arXiv:2412.06394.
 - Michal Kosinski. Evaluating large language models in theory of mind tasks. https://arxiv.org/abs/2302.02083, 2023.
 - Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tülu 3: Pushing frontiers in open language model post-training, 2024. URL https://arxiv.org/abs/2411.15124. arXiv preprint arXiv:2411.15124.
 - Adam Lerer, Natasha Gross, and Rob Fergus. Improving policies via search in cooperative partially observable games. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5847–5856. PMLR, 2020.
 - Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhai Wu, Behnam Neyshabur, Guy Gur-Ari, and Misra Vedant. Solving quantitative reasoning problems with language models. https://arxiv.org/abs/2206.14858, 2022. NeurIPS 2022 paper, arXiv:2206.14858.
 - Jiaqi Liu, Peng Hang, Xiaoxiang Na, Chao Huang, and Jian Sun. Cooperative decision-making for cavs at unsignalized intersections: A marl approach with attention and hierarchical game priors. *IEEE Transactions on Intelligent Transportation Systems*, 26(1):443–456, 2025.
 - Thang Luong and Edward Lockhart. Advanced version of gemini with deep think officially achieves gold-medal standard at the international mathematical olympiad, 2025. URL https://deepmind.google/discover/blog/advanced-version-of-gemini-with-deep-think-officially-achieves-gold-medal-standard
 - Hao Ma, Tianyi Hu, Zhiqiang Pu, Boyin Liu, Xiaolin Ai, Yanyan Liang, and Min Chen. Coevolving with the other you: Fine-tuning llm with sequential cooperative multi-agent reinforcement learning. *Neural Information Processing Systems (NIPS)*, 2024. doi: 10.48550/arXiv.2410.06101. URL https://arxiv.org/abs/2410.06101.
 - Meta AI. The llama 4 herd: The beginning of a new era of natively multimodal models. https://ai.meta.com/blog/llama-4-multimodal-intelligence/, 2025.
 - Mistral AI. Medium is the new large. https://mistral.ai/news/mistral-medium-3, May 2025.
 - Chunjiang Mu, Hao Guo, Yang Chen, Chen Shen, Die Hu, Shuyue Hu, and Zhen Wang. Multi-agent, human–agent and beyond: a survey on cooperation in social dilemmas. *Neurocomputing*, 610: 128514, 2024.

- OpenAI. GPT-4o System Card. arXiv preprint https://arxiv.org/abs/2410.21276,
- OpenAI. Competition shows humans are still better than ai at coding just. *The Guardian*, 2025a. URL https://www.theguardian.com/technology/2025/jul/26/competition-shows-humans-are-still-better-than-ai-at-coding-just. Published 26 July 2025.
 - OpenAI. Artificial intelligencer: Why ai's math gold wins matter. Reuters, 2025b. URL https://www.reuters.com/technology/artificial-intelligencer-why-ai-won-math-gold-2025-07-24/.
 - OpenAI. Introducing gpt-4.1 in the api. https://openai.com/index/gpt-4-1/, 2025a.
 - OpenAI. OpenAI o3 and o4-mini System Card. https://cdn.openai.com/pdf/ 2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card. pdf, 2025b.
 - David Premack and Guy Woodruff. Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences*, 1(4):515–526, 1978. doi: 10.1017/S0140525X00076512.
 - Qwen Team. Qwen3: Think deeper, act faster. https://qwenlm.github.io/blog/qwen3/, April 2025.
 - Aron Sarmasi, Timothy Zhang, Chu-Hung Cheng, Huyen Pham, Xuanchen Zhou, Duong Nguyen, Soumil Shekdar, and Joshua McCoy. Hoad: A dataset of hanabi self-play trajectories. https://ifaamas.org/Proceedings/aamas2021/pdfs/p1646.pdf, 2021.
 - Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback. *Neural Information Processing Systems (NIPS)*, 2020.
 - Arjun V. Sudhakar, Hadi Nekoei, Mathieu Reymond, Miao Liu, Janarthanan Rajendran, and Sarath Chandar. A generalist hanabi agent. https://arxiv.org/abs/2503.14555, 2025.
 - Haochen Sun, Shuwen Zhang, Lujie Niu, Lei Ren, Hao Xu, Hao Fu, Fangkun Zhao, Caixia Yuan, and Xiaojie Wang. Collab-overcooked: Benchmarking and evaluating large language models as collaborative agents. *arXiv preprint arXiv:2502.20073*, 2025.
 - Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances large language model capabilities. In *International Conference on Learning Representations* (*ICLR*), 2025.
 - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *CoRR*, abs/2201.11903, 2022. URL https://arxiv.org/abs/2201.11903.
 - Henry M. Wellman. The Child's Theory of Mind. Learning, Development, and Conceptual Change. MIT Press, Cambridge, MA, 1990. ISBN 9780262285940. doi: 10.7551/mitpress/1811.001.0001.
 - Zengqing Wu, Run Peng, Shuyuan Zheng, Qianying Liu, Xu Han, Brian I. Kwon, Makoto Onizuka, Shaojie Tang, and Chuan Xiao. Shall we team up: Exploring spontaneous cooperation of competing LLM agents. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 5163–5186. Association for Computational Linguistics, November 2024. doi: 10.18653/v1/2024.findings-emnlp.297. URL https://aclanthology.org/2024.findings-emnlp.297/.
 - xAI. Grok 3 beta the age of reasoning agents. https://x.ai/news/grok-3, February 2025.
 - Siwei Yang, Bingchen Zhao, and Cihang Xie. Aqa-bench: An interactive benchmark for evaluating llms' sequential reasoning ability. *Neural Information Processing Systems*, 2024. URL https://arxiv.org/abs/2402.09404.

C	ONTENTS	
1	Introduction	
2	Related Work	
3	Experiment Setup	
	3.1 Watson Setting	
	3.2 Sherlock Setting	
4	Benchmark Results	
5	Mycroft: Implicit Deductions from Multi-turn Play	
6	Ablations	
	6.1 Cross-Play	
	6.2 Instruction Tuning on HanabiLogs	
7	Future Work	
8	Conclusion	
9	Reproducibility Statement	
A	Why Hanabi?	1
В	LLM Agent Evaluation Suite	1
C	Model Analysis:	1
D	Hanabi Scores	1
E	Ablations	2
	E.1 Best-of-K Sampling	2
	E.2 Mixture of Agents	2
F	Human performance in Hanabi:	2
G	Seeds Evaluated	2
Н	Finetuning	2
	H.1 Training setup	2
	H.2 Results:	2
	H.3 Qualitative Example of Qwen behavior Change:	2
T	Single Agent Prompts	2

	I.1	MinCon Setup					
		I.1.1 Example Input and Output from OpenAI o4 mini:	26				
	I.2	DeductCon Setup:					
		I.2.1 Example input and output from OpenAI o4-mini:	28				
J	Best	of K - Final agent's Prompt (both MinCon and DeductCon setup)	32				
K	Mix	Mixture-of-Agent architecture					
L	L Example of MinCon Setup Multi-agent Prompts: L.1 Shared Information:						
M	Exa	mple of DeductCon Setup Multi-agent Prompts:	38				
N	N Multi-turn						
	N.1	Multi-turn Prompts	46				
	N.2	Multi-turn Results:	54				

A WHY HANABI?

Hanabi is a cooperative card game that has gained notable attention in the artificial intelligence research community as a benchmark for multi-agent coordination and reasoning under uncertainty Bauza (2010); Bard et al. (2020). The game involves 2-5 players working together to build firework displays by playing cards in ascending numerical order (1-5) across five different colors (red, yellow, green, blue, white). The fundamental challenge of Hanabi lies in its unique information structure: players can observe all cards held by their teammates but cannot see their own cards, creating an asymmetric information environment where successful play requires reasoning about what others know and communication through limited channels.

Players have access to a finite number of clue tokens (8 initially) that can be used to provide information about teammates' cards, indicating either all cards of a color or all cards of a rank in another player's hand. Additional clue tokens can be gained by discarding cards, but the maximum is capped at 8 tokens. This creates a tension between information gathering and resource management. The game's cooperative nature means all players share the same objective: maximize the collective score by successfully playing cards in the correct sequence while minimizing penalties from incorrect plays. The score is calculated as the sum of the highest card played in each color (e.g., if red reaches 4, blue reaches 3, green reaches 5, yellow reaches 2, and white reaches 1, the total score is 4+3+5+2+1=15). The maximum possible score is 25 (five colors \times five cards each), achieved by successfully completing all five firework displays. Each incorrect play consumes one of three fuse/life tokens, and the game immediately ends if all life tokens are exhausted. The game also ends when the deck becomes empty, after which players get one final round to play their remaining cards.

The shared objective, combined with information asymmetry, communication constraints, and the constant threat of game termination, creates a rich environment for studying collaborative decision-making and strategic reasoning. In Hanabi, all players must work toward a unified goal, collectively constructing ordered sequences of cards to maximize the team's score. This cooperative structure inherently differs from zero-sum or single-agent tasks, as success depends entirely on coordinated group performance rather than individual optimization. For LLMs, this means reasoning about collective utility functions and developing strategies that benefit the entire team, pushing models beyond self-interested decision-making paradigms. The game's core mechanism, where players observe others' cards but not their own creates a natural environment for testing theory of mind capabilities Premack & Woodruff (1978); Wellman (1990).

The variable player configurations in Hanabi introduce different strategic environments. While all games use the same 50-card deck, deck size and hand distributions vary: two and three-player games have 5 cards per hand (10 and 15 cards in hands, respectively), while four and five-player games use 4 cards per hand (16 and 20 cards in hands). The remaining deck size adjusts accordingly. These differences significantly impact the dynamics of cooperation. In two-player settings, direct one-to-one communication is sufficient. However, in other player settings, effective play requires distributed planning and multi-step coordination. For example, if player 4 needs to play a green 2 but cannot identify it, player 2 might give a rank clue ("2s"), and player 3 might then provide a color clue ("green"), allowing player 4 to deduce which of their card the green 2 is from the combined information. This interplay requires players to coordinate their clues and have a deep understanding of how each action advances the team's objective. This variety in configurations compels players to constantly consider their teammates' knowledge, beliefs, and potential deductions to make effective decisions. This mirrors the growing interest in assessing the theory of mind in large language models Kosinski (2023); Bubeck et al. (2023), while providing a more dynamic and impactful testing environment than traditional static psychological tasks.

An agent that performs consistently well across all player configurations demonstrates robust strategic understanding, rather than relying on brittle heuristics that overfit to specific scenarios. Because the optimal strategy differs drastically between player settings, consistent performance across them signals the development of generalizable reasoning principles. This cross-setting robustness is a crucial indicator of whether models have learned fundamental principles of cooperation and strategic reasoning, or simply developed configuration-specific patterns, making Hanabi an ideal benchmark for evaluating the generalizability of AI systems in varied collaborative environments.

B LLM AGENT EVALUATION SUITE

Our evaluation covered 17 LLMs across a spectrum of sizes, from 4B to over 600B parameters, spanning both open and closed-source families. We tested OpenAI models (o3, o4-mini OpenAI (2025b), GPT-4.1 GPT-4.1 mini OpenAI (2025a)), GPT-40 OpenAI (2024); Gemini (Gemini-2.5 Pro Comanici et al. (2025), Gemini-2.0 Flash Google DeepMind (2024), Gemini-2.5 Flash Google DeepMind (2025)); LLaMa-4 Maverick Meta AI (2025); DeepSeek-R1 (May 2025) DeepSeek-AI et al. (2025) and Deepseek-v3 (March 2025 DeepSeek-AI et al. (2024)); Qwen-3 (32B, 235B-A22B) Qwen Team (2025); Grok 3 and Grok 3-mini xAI (2025); Mistral 3 Medium Mistral AI (2025); and Claude Sonnet 3.7 Non-Thinking Anthropic (2025).

C MODEL ANALYSIS:

To better understand model performance, we analyzed the average number of turns played across 80 games (40 with the MinCon prompt, 40 with the DeductCon prompt), as shown in figure 8. Here, a "turn" denotes each instance the LLM was called during a game, summed across all players. Mistral Medium 3 and Llama Maverick typically failed early, averaging only about 20–25 turns per game, while most other models averaged over 60 turns in the MinCon prompt condition. In the DeductCon prompt scenario, most non-reasoning models (except GPT-4.1 and GPT-4.1 mini) quickly lost all three life tokens. Interestingly, there was no direct correlation between the number of turns played and final scores: top-performing models played slightly fewer turns than others such as GPT-4.1 and GPT-4.1 mini. This suggests that stronger reasoning models were more efficient in maximizing rewards per turn. In general, all models played fewer turns with the DeductCon prompt, except for Mistral Medium 3. For reasoning models, prompt type had little effect on turns played, aside from cases like Qwen-235B-A22B, which sometimes lost life tokens faster and ended games earlier with the DeductCon prompt. In contrast, non-reasoning models, except for the GPT-4.1 family, played significantly fewer turns with the DeductCon prompt, suggesting they often failed by losing all life tokens earlier compared to the MinCon prompt.

We further investigated why non-reasoning models struggled in the DeductCon prompt case. When given simple, rigid prompts such as "always play the safe move," non-reasoning models generally succeeded. However, with more complex instructions that required probability calculation, these models often became confused. In contrast, reasoning models handled multiple objectives well,

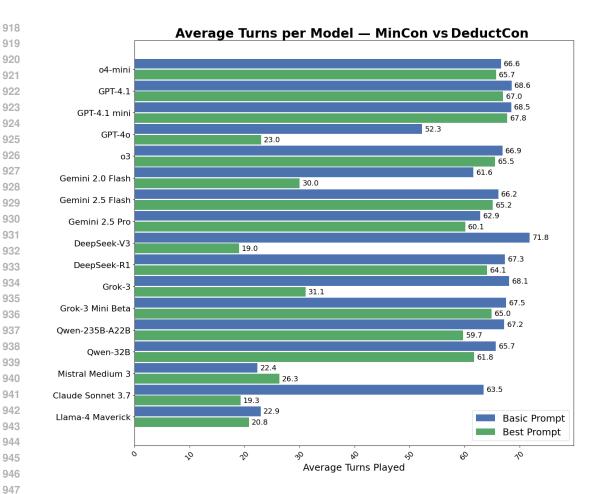


Figure 8: Average number of turns played by each model, averaged over the two- through five-player settings.

including calculating probabilities, providing reasoning, and following instructions to output in the desired JSON format.

Non-reasoning models like Llama 4 Maverick frequently made high-risk plays without sufficient information, leading to rapid loss of life tokens and early game termination. Gemini 2.0 Flash was more cautious in the MinCon prompt scenario but often gave redundant clues and made unnecessary discards, resulting in lower scores despite playing approximately three times more turns than Llama 4 Maverick. GPT-40 showed significant inefficiencies as well, frequently giving repetitive clues and misplaying by failing to track the game state, which hurt its overall performance even with a high number of turns. Mistral Medium 3 tended to prioritize giving information over executing clear plays; once out of information tokens, it would play or discard cards at random, making it the weakest performer in this group. However, its performance improved considerably when given more contextual information, highlighting that it lacked world knowledge about Hanabi.

We also observed several peculiar behaviors. Models sometimes assigned higher ratings to moves they did not select. This behavior was more common in non-reasoning models than in reasoning models. Some models attempted to play higher-numbered cards onto fireworks stacks that had not yet reached the required lower numbers, resulting in life token loss. For example, when the green firework was at 2, the model played a green 5, justifying the move by claiming it would increase the score by three. This occurred despite explicit instructions in the prompt that fireworks must be built sequentially. Each model family posed distinct challenges: for example, GPT-40 occasionally output invalid moves; Qwen, DeepSeek, and Gemini family models sometimes failed to follow instructions, producing outputs in an incorrect format and causing experiment failures. Because Hanabi is a

sequential game, such inconsistencies necessitate robust code capable of either repeatedly recalling the API until a valid result is obtained, or if repeated attempts fail parsing all prior valid moves and resuming play from that point. We advise future work with the Hanabi Learning Environment to anticipate and accommodate these issues.

D HANABI SCORES

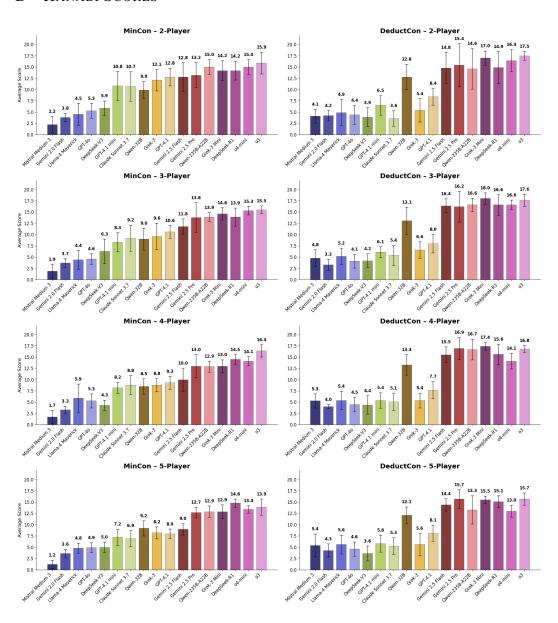


Figure 9: Performance of various LLMs on the Hanabi benchmark across two- to five-player settings. The left column shows average scores (over 10 seeds) of the MinCon Prompt, right column shows the average scores of the DeductCon Prompts.

E ABLATIONS

A single Hanabi game typically requires at least 60 turns (Figure 8). Due to the non-deterministic nature of LLM outputs, the quality of reasoning can vary across runs. We examine this behavior empirically with Best-of-K sampling (Section E.1) and a Mixture of Agents approach (Section E.2).

E.1 BEST-OF-K SAMPLING

To improve reliability, we use Best-of-K sampling Stiennon et al. (2020): for each turn, we sample the agent k times, generating multiple candidate actions (which may not all be unique), and then prompt the agent to select the single best option from these samples. See Appendix J for details of the prompts used. For our Best-of-K experiments, similar to our prompting strategy ablations (Section 3.2) we used Grok-3-mini in the 5-player setting with a fixed seed (3), running each configuration 10 times.

Varying K. We evaluate performance for k = 1, 2, 3, 4, 5, 6, and 7, with the MinCon prompt, SPIN-Bench prompt, and our DeductCon prompt, where each agent is given the same prompt k times. As shown in Figure 10, for k = 1 and 2 our DeductCon prompt outperforms the others, as previously discussed in Section 3.2. However, as k increases, our DeductCon prompt performance converges with SPIN-Bench. While baselines improve until k = 5 and then dip, our DeductCon prompt shows consistent performance across all k values (sample variance $\sigma = 1.23$ on 0 to 25 scoring scale), with minimal gains from increased sampling. There is also a clear performance gap (> 1.5 on average across K values) between the MinCon prompt and the other two setups.

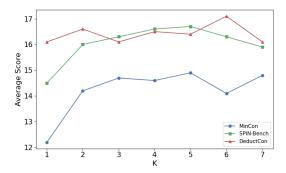


Figure 10: Best-of-K average Hanabi score with the MinCon prompt, SPIN-Bench prompt, and our DeductCon prompt, averaged over 10 runs on the 5-player Seed 3 setting.

Varying # Players. To compare Best-of-K performance across player counts ($2\ to\ 5$) and context (MinCon and DeductCon prompts), we fix k=5, as for both SPIN-Bench and MinCon prompt setups, this is where game scores peak (Figure 10). We find that our DeductCon prompt consistently outperforms the MinCon prompt across all player counts with Best-of-5 sampling, which we show in Figure 11. We also compare Best-of-5 sampling to Best-of-1 (i.e. K=1, no sampling), which we have already shown in Figure 4. We observe that for Grok-3-mini, using Best-of-5 sampling with the MinCon prompt improves performance over K=1 in all cases (+1.5 on average) except the 2-player setting (-0.1). In contrast, applying Best-of-5 to the DeductCon prompt across 40 games yields negligible further improvement (+0.1 on average) compared to K=1, which is consistent with our observations while varying K in Figure 10.

E.2 MIXTURE OF AGENTS

With our DeductCon prompt, we observed that sampling from K agents using the *same* prompt gave no score benefits as agents would often select consistent actions even as K increased. To encourage diversity in agent selected actions, inspired by Mixture of Agents (MoA), Wang et al. (2025) we use five parallel agents with specific roles to generate diverse outputs, which are then provided to an aggregator agent for final move selection. As prior work Wei et al. (2022) and our single-agent experiments (Section 4) demonstrated that better prefill improves agent performance, we ensured that all parallel agents supplied detailed, relevant, and diverse information to the final agent. See

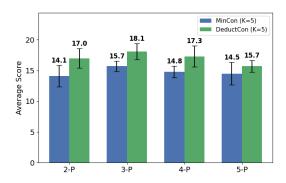


Figure 11: Best-of-K average Hanabi score at K=5, comparing the Watson and Sherlock prompts across player count (2-5).

Appendix L for MinCon and DeductCon multi-agent prompting details, as well as rubrics used by some of the agents below:

Agent 1 (MinCon): In both setups, this agent used the same prompt as the single-agent baseline. **Agent 2 (Clue Preference):** Same prompt as Agent 1 with an additional instruction to choose rank clues over color clues when both were equally favorable.

Agent 3 (Analyst): Required to provide analysis for all cards in the agent's and other players' hands. In the MinCon prompt, we observed that the aggregator agent often based its answer on the Analyst's response. Therefore, in the DeductCon prompt, we asked the agent to follow a detailed rubric which provided comprehensive information for each card.

Agent 4 (Discard): Tasked with identifying safe and critical discards. The DeductCon prompt uses a rubric for more structured prefill to the aggregator agent.

Agent 5 (History): This agent infer teammates' intentions based on prior move history (10 moves for the MinCon prompt, full history for the DeductCon prompt). We observed that with MinCon, this agent contributed only generic information that the aggregator ignored. With DeductCon, we included in-context examples to encourage the agent to speculate more actively.

Agent 6 (Aggregator): Receives all specialist agent outputs along with the game state and history to select the mixture of agents' final move. See Appendix K for a detailed setup of our mixture of agents and Appendix M for all the prompts.

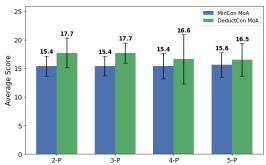


Figure 12: Mixture of Agents (MoA) average score with the MinCon and DeductCon prompting strategies across 2, 3, 4, and 5-player settings. All player count settings use the six agents described in Section E.2, except 2-player, which omits the History Agent.

With our mixture of agents framework, as shown in Figure 12, we observed that 5-player score improves with both MinCon (+1.1) and DeductCon (+0.8) settings compared to Best-of-5 sampling. Mixture of agent scores are similar to Best-of-5 for the 3-player and 4-player games (+0.3 for MinCon and -0.5 for DeductCon). With the DeductCon prompt, in 4 and 5 player settings, one run ended prematurely, which lowered the overall mean and increased the standard deviation. Omitting this outlier run results in 4-player score 17.89 (+0.6 over Best-of-K) and 5-player score 17.34 (+1.6 over Best-of-K). High score variance was most pronounced in the 2-player setting: the

history agent's speculation led to highly variable results (with one run scoring 23, while a few others scored below 10). As a result, we removed the history agent for the 2-player setting.

Takeaways. We find that reasoning models excel at following explicit instructions and perform at the third quartile (75th percentile) of human players from BoardGameGeek (see Appendix F). However, they often fail to anticipate the likely actions of other players. To reach the top 25th percentile, future models may need to be explicitly trained on theory of mind tasks. Our experiments with prefilled prompts (Figure 3) show that reasoning models rarely perform worse when provided with richer, relevant context and instruction (in our case, the Sherlock prompt). This suggests that further improvements are possible if agents are exposed to more in-context strategy specific to different player settings alongside additional Hanabi domain knowledge.

F HUMAN PERFORMANCE IN HANABI:

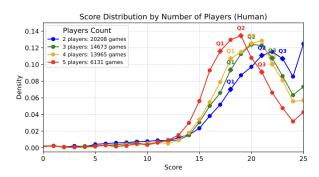


Figure 13: Distribution of human Hanabi scores (2–5 players) collected from BoardGameGeek. The graph is taken from SPIN-Bench Yao et al. (2024).

We use the human baseline provided by SPIN-Bench Yao et al. (2024), which aggregated 54,977 human-played Hanabi games from BoardGameGeek, covering 2- to 5-player settings. Our reasoning models reach the Q1 threshold in self-play, indicating they now perform comparably to the lower quartile of human players, but still lag behind the median (Q2) and upper quartile (Q3) benchmarks.

G SEEDS EVALUATED

All the models were evaluated on seeds 1,2,3,5,7,11,13,17,19,23 except for the Qwen-3-4B-Instruct-2507 models to avoid memorization effects. We evaluated them on Seeds 4,6,8,10,12.

H FINETUNING

H.1 TRAINING SETUP

Data. We fine-tune on **HanabiLogs** (ours), formatting each record with the model's chat template and applying response-only supervision (tokens before the assistant span labeled -100), while restricting the corpus to outputs from grok3mini or o3 Sherlock setup outputs.

Main model. We train <code>Qwen/Qwen3-4B-Instruct-2507</code> with LoRA $(r=16, \alpha=32, dropout=0.05)$ on attention and MLP projections, using AdamW, bf16, gradient checkpointing, and sequence chunking with <code>block_size</code> and <code>doc_stride</code>. Unless noted: <code>lr=2e-5</code>, <code>per_device_batch_size=2</code>, <code>grad_accum=8</code>, <code>num_train_epochs=3</code>, <code>block_size=16384</code>, <code>doc_stride=256</code>.

H.2 RESULTS:

Instruction tuning/SFT significantly improved model performance in the 2- and 3-player settings compared to the 4- and 5-player settings because the model learnt basic strategies, such as playing

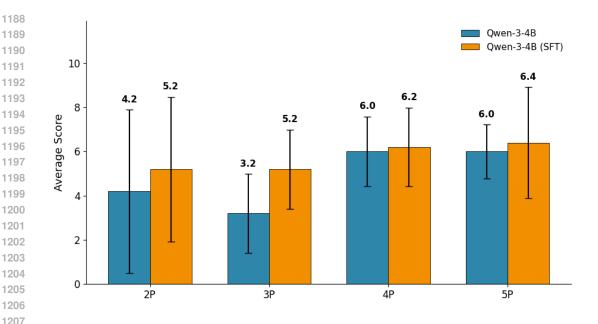


Figure 14: Average scores of Qwen-3-4B-Instruct-2507 before and after SFT across 2-5 player settings.

rank 1 initially and taking risks at the final turn. This also made the model overconfident at times, which resulted in early exit in 4 and 5-player settings.

H.3 QUALITATIVE EXAMPLE OF QWEN BEHAVIOR CHANGE:

We illustrate a behavioral shift in <code>Qwen</code> model after supervised fine-tuning (SFT). Before fine-tuning, the models did not apply the opening heuristic: when all firework stacks are at 0, any card known to be rank 1 is safe and will increase the score by 1. After SFT, the models consistently adopt this strategy.

Game State

1208

1209

1210 1211

1212

1213 1214 1215

1216

1217

1218

1219

1220

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1236

1237

1239

1240

1241

There are 3 life tokens and 8 info tokens remaining. The fireworks progress: R stack is at 0, Y stack is at 0, G stack is at 0, W stack is at 0, B stack is at 0. Your hand contains the following cards: Card 0:- Known info: 'XX'. No hints about this card's color or rank have been given yet.- Could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5.Card 1:- Known info: 'XX'. No hints about this card's color or rank have been given yet.- Could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5.Card 2:- Known info: 'XX'. No hints about this card's color or rank have been given yet.- Could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5.Card 3:-Known info: 'XX'. No hints about this card's color or rank have been given yet.- Could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5.Card 4:- Known info: 'XX'. No hints about this card's color or rank have been given yet.- Could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5. From your perspective, you can see the other players' hands clearly. Here's what youobserve:Player +1's hand:- A card: You can see the card: 'R3', This player has no specific hints about the card'sidentity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5.- A card: You can see the card: 'Y1', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5.- A card: You can see the card: 'W3', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5. - A card: You can see the card: 'B2', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5. - A card: You can see the card: 'W1', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5. There are 40 cards remaining in the deck. The discard pile contains: no cards discarded yet.

1243 1244

1245

1246

1247

1248

1249

1250

1251

1253

1255

1257

1259

1261

1262

1263

1264

1265

1266

1267

1268

1270

1276

1278

1279

1280

1281

1282

1283

1284

1285

1286

1291

1293

1295

"reason": "The current state shows no firework stacks have been built (all at 0), and no cards have been played or discarded. All players' hands are visible, and we have 3 life tokens and 8 info tokens. My hand contains 5 unknown cards, each with equal probability across all colors and ranks. However, I can observe that Player +1 holds several cards: R3, Y1, W3, B2, W1. These provide critical information about the distribution of ranks and colors. Key observations: - Player +1 has a Y1, which means yellow has at least one rank 1. This supports starting a yellow firework with rank 1. - Player +1 has R3, W3, B2 - suggesting that rank 3 appears in red, white, and blue, and rank 2 in blue. This implies that rank 3 is common, and rank 2 is present in blue. - Rank 5 is extremely rare (only one per color), so we should preserve any card that might be rank 5. - Rank 1 is common (3 per color), so it is safe to use for starting sequences. Since no firework stacks exist, we cannot play any card yet - playing a card without a valid sequence would be illegal and cost a life token. Therefore, any play action (0-4) would be risky and illegal unless we are certain of the sequence. Instead, we should use a hint to help a teammate identify a card they can play. Player +1 has Y1, which is a valid starting card for yellow. However, they may not know that yellow has a rank 1. If they play Y1, it would start the yellow firework. But they might not realize it is playable. We can give a hint about the color yellow to Player +1 to confirm that yellow has a rank 1. This is a safe, low-risk hint that directly supports a valid play. Since Player +1 already has Y1, they know it exists, but they may not know it is playable or that it is the correct rank to start. However, the hint options are only about revealing a color or rank to another player. We can give a hint about a color (e.g., Reveal player +1 color Y) to tell them that yellow is present in their hand. But they already know that from their own hand. Alternatively, we can hint about a rank. For example, if we hint that Player +1 has rank 1, that would confirm they have at least one rank 1 card, which they already do (Y1). So that is redundant. But if we hint about a color that they do not have, that would be false. We must only give accurate hints. Player +1 has: R3, Y1, W3, B2, W1 -> so they have: - Red: R3 - Yellow: Y1 - White: W3, W1 - Blue: B2. So they have no green cards. That is critical. Therefore, we can give a hint about green to Player +1 – but only if we know they do not have it. Since they do not have green, we cannot hint that they have green. We can only give accurate hints. We can only say they have a color or rank that they actually have. So the only valid hints are about colors they already have or ranks they already have. But we can give a hint about a rank they have – for example, rank 1 – to confirm they have a rank 1 card. They do (Y1 and W1). So we could hint Reveal player +1 rank 1. This would tell them that they have at least one rank 1 card, which they already know. So it is redundant. Similarly, they have rank 2 (B2), rank 3 (R3, W3), so we could hint rank 2 or 3. But giving a hint about a rank they already have does not add value. Therefore, no useful hint can be given to Player +1 about their own hand, because all cards are visible. But the game allows hints to other players about colors or ranks they have. So if we say Reveal player +1 color Y, that means we are telling them that yellow is present in their hand, which is true. But they already know that from their own hand. So all hints about colors or ranks they already have are redundant. Therefore, no hint is useful at this stage. So what should we do? We must either play a card or discard. But playing any card now would be illegal, because no firework sequence exists. We cannot play a card unless it matches the current top of a stack (which is 0 for all colors). So we cannot play any card. Therefore, all play actions (0-4) are illegal and would cost a life token. Thus, we must discard a card to gain an info token. We should discard a card that is least likely to be useful, ideally one that is not rank 1, not rank 2,3,4, and not a critical color. But we have no information about the cards' colors or ranks. We must use probability. Total deck: 50 cards. Cards per color: - Rank 1: 3 per color -> 15 total - Rank 2: 2 per color -> 10 total - Rank 3: 2 per color -> 10 total - Rank 4: 2 per color -> 10 total - Rank 5: 1 per color -> 5 total. Cards visible: Player +1 has: - R3 (red 3) -Y1 (yellow 1) - W3 (white 3) - B2 (blue 2) - W1 (white 1). So removed from deck: - Red: one 3 - Yellow: one 1 - White: one 3 and one 1 - Blue: one 2. So remaining deck: 45 cards. Now, in our hand, each card has equal probability across all colors and ranks, but we must adjust for known cards. For example, red:

1297

1298

1299

1300

1301

1302

1303

1304

1305

1309

1310

1311

1312

1315

1316

1317

1318 1319 1320

1321 1322

1328

1330

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1347 1348

rank 3 is taken -> so red 3 is gone, but red 1,2,4,5 still exist. Yellow: rank 1 is taken -> so yellow 1 gone, but 2,3,4,5 still exist. White: ranks 1 and 3 taken -> so white 1,3 gone -> 2,4,5 remain. Blue: rank 2 taken -> 1,3,4,5 remain. Green: no cards seen -> all ranks 1-5 still available. So for a card in our hand: it could be green, red, yellow, white, or blue. Green has full distribution. Red is missing 3 (possible 1,2,4,5). Yellow is missing 1 (possible 2,3,4,5). White is missing 1 and 3 (possible 2,4,5). Blue is missing 2 (possible 1,3,4,5). Now, which card to discard? We want to discard a card that is least likely to be useful. Critical cards to preserve: - Rank 5 (only one per color), very valuable for completing sequences. - Ranks 2,3,4, especially 3 and 4, which are common. - Color green, not yet seen, so may be key. So we should avoid discarding a card that might be a 5 (in any color). But we have no way to know. So we discard a card that is least likely to be 5 or 3 or 4. But all ranks are similarly likely given our lack of info. Alternatively, we can discard a card that is not in a color that has a high probability of being used, but all colors are symmetric. Best strategy without info: discard a card that is not rank 1, because 1s are common and useful for starting sequences, yet we do not know the rank. So we can only discard based on probability. Probabilities (given 45 unknown cards): P(5)=5/45=0.111..., P(1)=13/45=0.289..., P(2)=9/45=0.2, P(3)=8/45=0.178..., P(4)=10/45=0.222.... So 5s are rare but critical; discarding at random risks losing one, but there is no better choice without info. Initially I pick discarding card index 2 as an arbitrary choice. However, checking the provided legal actions shows only plays (0-4) and reveals (5-11); there is no discard option in this mapping. Therefore discard is not available right now. Given that all plays would be illegal (no stacks started), the only remaining legal move type is a hint. Among available hints to Player +1 (colors R/Y/W/B or ranks 1/2/3), all are somewhat redundant since they already hold those. Still, we must act, and a truthful hint is preferable to an illegal play. I choose to hint rank 3 to Player +1. Action: 11 -> Reveal player +1 rank 3. This is safe, accurate, and provides minimal but useful information. Rating: 0.0, because it is redundant, but better than nothing. Earlier I thought play action 2 had rating 0.9 because I thought card 2 might be playable, but it is not. So all play actions are illegal. Therefore, no play action is safe. Thus, the only safe action is to give a hint. I choose to give a hint about rank 3 to Player +1. Final decision: action 11.", "action": 11

```
Owen-3-4B-Instruct-2507 Response - turn 1
```

"reason": "The current game state shows all fireworks at 0, meaning only rank 1 cards can be played safely. Player +1 has a Y1 card that is playable, but they lack information to identify it. Hinting rank 1 to Player +1 will reveal exactly one card (Y1) and allow them to play it on their next turn, scoring a point with no risk. This is prioritized over other hints or plays because it directly enables a safe play, while my own hand has no playable cards due to unknown identities and low probability of success (only about 1/5 chance for any card to be rank 1). Discarding would gain an info token but delays scoring, and other hints are less precise or less immediately useful. With 8 info tokens available, spending one is cost-effective.",

```
"action": 9
```

From above, we can notice that the models response become shorter due to training on outputs of o3 and grok3mini. In the next turn, the other player played the rank 1 card.

```
1350
                   {"action": 1, "rating": -0.5},
1351
                   {"action": 2, "rating": 0.2},
1352
                   {"action": 3, "rating": 0.2},
1353
                   {"action": 4, "rating": -0.5},
1354
                   {"action": 5, "rating": -0.8},
                   {"action": 6, "rating": 0.9}, 
{"action": 7, "rating": -0.8},
1355
1356
                   {"action": 8, "rating": -0.8},
1357
                   {"action": 9, "rating": 0.9},
1358
                   {"action": 10, "rating": 0.1},
1359
                   {"action": 11, "rating": 0.1},
                   {"action": 12, "rating": 0.1},
1360
                   {"action": 13, "rating": 0.1},
                   {"action": 14, "rating": 0.1},
                   {"action": 15, "rating": 0.1},
1363
                   {"action": 16, "rating": 0.1}
1364
1365
```

"reason": "We have 7 information tokens, so we can afford to spend one on a clue later. The safest and highest-value move right now is to play one of our known rank-1 cards. - Card 1 and Card 4 are both rank-1, and every color is currently at 0, so either one is guaranteed to succeed and advance the fireworks. - Playing a rank-1 card costs no information tokens, so we keep our supply at 7. - Discarding a rank-1 card would give us an info token but risks throwing away a critical card (we still have 2 rank-1s left in the deck, so we can afford to lose one, but it is not as good as scoring a point). - Playing an unknown card (0, 2, or 3) is risky and could cost a life. - Giving a clue now would delay scoring and might not lead to an immediate play (e.g., the other player might not know which rank-1 to play). Therefore, playing a rank-1 card maximizes our score with minimal risk. I will choose to play Card 1 (the first rank-1 in my hand).",

"action": 6

1367

1369

1370

1371

1372

1377

1379 1380

1381

1382

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1399

1400

1401

1402

1403

I SINGLE AGENT PROMPTS

I.1 MINCON SETUP

I.1.1 EXAMPLE INPUT AND OUTPUT FROM OPENAI O4 MINI:

MinCon Prompt & Response

Input Prompt:

{ "system_prompt": "You are an expert AI in the cooperative card game Hanabi. Your goal is to help the team achieve the highest possible score (max 25).

Analyze the entire game state provided, including your hand knowledge, visible hands of other players, fireworks, discards, deck size, lives, and info tokens.

Consider all strategic priorities:

- 1. Safe Plays: Prioritize playing cards you KNOW are playable on the fireworks.
- 2. **Useful Clues:** If no safe play and info tokens > 0, consider giving clues that enable immediate plays, save critical cards, or provide significant new information without being redundant.
- 3. **Safe Discards:** If no safe play and no high-value clue (or info tokens == 0), discard the safest possible card.

Explain your reasoning clearly, referencing the game state, and then state your chosen move number.

OUTPUT FORMAT:

Reasoning: [Your detailed reasoning justifying your choice based on the game state and strategic priorities]

Chosen Move Number: [number]",

"user_prompt": "You are Player 1. Analyze the game state and propose the best move number. Explain your reasoning.

```
1404
          Game State:
1405
          P1 (2p Game). Lives: 3, Info: 1, Deck: 0.
1406
          Fireworks: R2 Y4 G2 W3 B1. Discards: 1 red card rank 1, 1 red card rank 2, 2 red cards rank 3,
1407
          1 red card rank 4, 2 yellow cards rank 1, 1 yellow card rank 2, 1 yellow card rank 4, 1 yellow
1408
          card rank 5, 1 green card rank 1, 1 green card rank 2, 1 green card rank 3, 2 green cards rank 4, 1
1409
          green card rank 5, 2 white cards rank 1, 1 white card rank 2, 1 white card rank 3, 2 white cards
1410
          rank 4, 2 blue cards rank 1, 1 blue card rank 2, 2 blue cards rank 4, 1 blue card rank 5.
1411
          FINAL ROUND! 1 turns left (P1 is last).
1412
          Visible Hands: P0 Hand: [Blue 3, Red 1, Red 4, White 5, Blue 3]
1413
          Your Knowledge (Hints): [UnknownColor UnknownRank, UnknownColor UnknownRank,
1414
          UnknownColor 3, Blue 2, UnknownColor UnknownRank] (Indices 0–4).
          Others' Knowledge: P0 Knows: [UnknownColor 3, UnknownColor 1, UnknownColor
1415
          UnknownRank, UnknownColor UnknownRank, UnknownColor UnknownRank]
1416
1417
          Legal Moves:
1418
          0. (Discard 0)
1419
          1. (Discard 1)
1420
          2. (Discard 2)
1421
          3. (Discard 3)
1422
          4. (Discard 4)
1423
          5. (Play 0)
1424
          6. (Play 1)
1425
          7. (Play 2)
          8. (Play 3)
1426
          9. (Play 4)
1427
          10. (Reveal player +1 color R)
1428
          11. (Reveal player +1 color W)
1429
          12. (Reveal player +1 color B)
1430
          13. (Reveal player +1 rank 1)
1431
          14. (Reveal player +1 rank 3)
1432
          15. (Reveal player +1 rank 4)
1433
          16. (Reveal player +1 rank 5)
1434
1435
          Current Info Tokens: 1
1436
          Output Format:
1437
          Reasoning: [Your detailed reasoning justifying your choice based on the game state and strategic
1438
          priorities]
1439
          Move Ratings: [Rate each legal move from -1 (terrible) to 1 (excellent), like "Move 0: 0.5, Move
1440
          1: -0.3, Move 2: 1.0, ..."]
1441
          Chosen Move Number: [number]"
1442
1443
           response: "Reasoning:
1444
          It's the last turn (no cards left in deck, P1's only action). The current fireworks are R2, Y4, G2,
1445
          W3, B1 (total 12). The only way to increase score is to play a card now. You know your card
1446
          at index 3 is Blue 2, which is exactly the next playable card on the blue firework (B1\rightarrowB2).
1447
          Playing it gains +1 point, raising the score to 13.
1448
          All other play slots are uncertain and risk a misplay for no guaranteed score; discards or clues
1449
          yield zero benefit (P0 won't get another action and you can't draw). Therefore the uniquely
1450
          correct move is to play card 3.
1451
1452
          Move Ratings:
1453
          Move 0 (Discard 0): -0.5
          Move 1 (Discard 1): -0.5
1454
          Move 2 (Discard 2): -0.5
1455
          Move 3 (Discard 3): -0.5
1456
          Move 4 (Discard 4): -0.5
1457
```

```
1458
          Move 5 (Play 0): 0.0
1459
          Move 6 (Play 1): 0.0
1460
          Move 7 (Play 2): 0.2
1461
          Move 8 (Play 3): 1.0
1462
          Move 9 (Play 4): 0.0
1463
          Move 10 (Clue P0 color R): -1.0
1464
          Move 11 (Clue P0 color W): -1.0
1465
          Move 12 (Clue P0 color B): -1.0
1466
          Move 13 (Clue P0 rank 1): -1.0
1467
          Move 14 (Clue P0 rank 3): -1.0
1468
          Move 15 (Clue P0 rank 4): -1.0
1469
          Move 16 (Clue P0 rank 5): -1.0
1470
          Chosen Move Number: 8"
1471
1472
```

Due to a prompt oversight, the system prompt included an output format without move ratings, while the user prompt specified move ratings. All models followed the user prompt as intended. As this was the lower bound case, we retained this setup.

I.2 DEDUCTCON SETUP:

1473

1474

1476

1477 1478

1479

1480 1481

1482

1483

1484

1485

1486

1487 1488

1489

1490 1491

1492

1493

1494

1495

1496

1497 1498

1499

1500

1501

1502

1503

1506

1507

1509

1510

1511

I.2.1 EXAMPLE INPUT AND OUTPUT FROM OPENAI 04-MINI:

DeductCon Prompt & Response

"user_prompt": "You are a master of hanabi game. You are playing a game of Hanabi with 2 players. Hanabi is a cooperative card game where players work together to create a series of fireworks by playing cards in ascending numerical order starting from 1. Each player holds their cards facing outward so that all players can see everyone else's cards but not their own. The objective is to play cards in sequence (1 through 5) for each color without making mistakes. There are 5 different colors and each color has cards numbered 1 to 5.

Key Rules:

On your turn, you have three types of possible actions:

Give a Hint(Reveal): Provide a hint to another player about their cards, specifying either a color or a number present in their hand. Hints must be accurate and can only reveal positions of cards matching the hint.

Discard a Card: Discard one of your own cards to potentially gain an Info token.

Play a Card: Attempt to play a card from your hand. If played correctly in sequence, it adds to the fireworks; if not, it reduces one life token.

Tokens:

Life Tokens: Deducted when a wrong card is played.

Info Tokens: Used to give clues.

Illegal Moves: Playing a card that cannot be placed properly costs a life token. If life tokens reach zero, the game ends in failure.

Game End: The game ends when all fireworks are completed (perfect score of 25), or when the deck is exhausted and each player has taken one final turn, or when the players run out of life tokens.

State Representation: The game state is represented with the following details:

Life tokens: Number of remaining life tokens.

Info tokens: Number of available information tokens.

Fireworks: Current progress on each firework color (e.g., R1, Y0, G1, W0, B0).

Discards: Cards that have been discarded.

1512 1513 Your Role: 1515 You are one of the players, cooperating with others to maximize the total score of the fireworks (the number of cards correctly played in sequence). 1516 Although you cannot see your own cards, you can see the cards in the hands of your teammates. 1517 Use hints, discards, and plays strategically to guide the team towards successful sequences. 1518 1519 Remember, communication is limited to hints about colors or numbers only, and shar-1520 ing illegal or extraneous information is not allowed. Work together, follow the rules, and aim for 1521 the highest cooperative score possible! 1522 1523 Below is the current detailed state information. 1524 1525 **Game State:** 1526 There are 3 life tokens and 2 info tokens remaining. The fireworks progress: R stack is at 5, Y stack is at 5, G stack is at 3, W stack is at 2, B stack is Your hand contains the following cards: 1529 Card 0: 1530 - Known info: 'XX'. No hints about this card's color or rank have been given yet. 1531 - Could be any of these colors: Red, Yellow, Green, White with ranks: 1, 3, 4, 5. 1532 1533 - Known info: 'XX'. No hints about this card's color or rank have been given yet. 1534 - Could be any of these colors: Red, Yellow, Green, White with ranks: 1, 2, 3, 4, 5. 1535 Card 2: 1536 - Known info: 'XX'. No hints about this card's color or rank have been given yet. 1537 - Could be any of these colors: Red, Yellow, Green, White with ranks: 1, 2, 3, 4, 5. 1538 - Known info: 'BX'. Known: color is blue. 1539 - Could be any of these colors: Blue with ranks: 1, 2, 3, 4, 5. 1540 Card 4: 1541 - Known info: 'XX'. No hints about this card's color or rank have been given yet. 1542 - Could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5. 1543 From your perspective, you can see the other players' hands clearly. Here's what you observe: Player +1's hand: 1545 - A card: You can see the card: 'W1', This player has no specific hints about the card's identity, 1546 This player knows it could be any of these colors: Yellow, Green, White with ranks: 1, 2, 3. 1547 - A card: You can see the card: 'W2', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Red, Yellow, Green, White with ranks: 1, 2, 3. - A card: You can see the card: 'Y4', This player has no specific hints about the card's identity, 1549 This player knows it could be any of these colors: Red, Yellow, Green, White with ranks: 1, 2, 3, 1550 1551 - A card: You can see the card: 'R3', This player has no specific hints about the card's identity, 1552 This player knows it could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1553 1, 2, 3, 4, 5. 1554 There are 0 cards remaining in the deck. The discard pile contains: 2 red cards rank 1, 1 red card 1555 rank 4, 1 yellow card rank 1, 1 yellow card rank 2, 1 yellow card rank 3, 2 green cards rank 1, 1 1556 green card rank 2, 1 green card rank 3, 2 green cards rank 4, 1 green card rank 5, 1 white card 1557 rank 1, 2 white cards rank 3, 1 white card rank 5, 2 blue cards rank 1, 1 blue card rank 2, 1 blue card rank 3, 1 blue card rank 5. FINAL ROUND: The deck is empty. You are the final player and this is the final turn 1560 for the whole game. 1561 Please think step by step based on the current state 1563

Think step by step

1566 ## Evaluate Playable Cards in Hand 1567 1568 Look at each card in your hand. 1569 Cross-reference with the current game state to see if any card can be immediately played to 1570 complete or extend a firework stack. 1571 Consider hints you have received about each card (color/rank information) to determine if it 1572 might be safe to play. 1573 If a card can be played without risk, prioritize playing it to score a point. 1574 1575 ## Consider Teammates' Hands and Hint Opportunities 1576 Analyze the visible cards in your teammates' hands. Identify if any of their cards can now be played based on the current firework stacks or previous 1579 If you notice a teammate holds a card that can be played but they may not realize it, think about 1580 what hints you could give them. 1581 Use hints to communicate critical information, such as color or rank, to help them make the right play. Choose the hint that maximizes the chance for a correct play while considering the limited hint 1584 tokens. 1585 1586 ## Assess Discard Options to Gain Info Tokens 1587 Look for cards in your hand that are least likely to be playable or helpful in the near 1588 Consider the remaining deck composition and cards already played/discarded to predict the 1590 value of each card. 1591 Discard a card that you believe to be least useful to gain an Info token, especially if no immediate 1592 playable or hint options are available. 1593 Ensure that discarding this card won't permanently remove a critical card needed to complete 1594 any firework stack. 1595 1596 Now it's your turn. You can choose from the following legal actions: 1597 The legal actions are provided in a mapping of action identifiers to their descriptions: {0: '((Discard 0))', 1: '((Discard 1))', 2: '((Discard 2))', 3: '((Discard 3))', 4: '((Discard 4))', 5: '((Play 0))', 6: '((Play 1))', 7: '((Play 2))', 8: '((Play 3))', 9: '((Play 4))', 10: '((Reveal player +1 color R))', 11: '((Reveal player +1 color Y))', 12: '((Reveal player +1 color W))', 13: '((Reveal player +1 rank 1))', 14: '((Reveal player +1 rank 2))', 15: '((Reveal player +1 rank 3))', 16: '((Reveal player +1 rank 4))'} 1604 (Reveal player +N color C): Give a hint about color C to the player who is N positions 1605 ahead of you. (Reveal player +N rank R): Give a hint about rank R to the player who is N positions ahead. (Play X): Play the card in position X from your hand (Card 0, Card 1, Card 2, etc.). 1608 (Discard X): Discard the card in position X from your hand (Card 0, Card 1, Card 2, etc.). 1609 Based on the annotated state and the list of legal actions, decide on the most appropri-1610 ate move to make. Consider factors like current tokens, firework progress, and information 1611 available in hands. Then, output one of the legal action descriptions as your chosen action. 1612 1613 Your output should be in this format: 1614 1615 "reason": string, 1616 "action": int} And the action should be one of the legal actions provided above. 1617 You can only use json valid characters. When you write json, all the elements (including all the 1618

keys and values) should be enclosed in double quotes!!!

```
1620
1621
           CRITICAL: Also include move ratings in this exact JSON format:
1622
1623
           "move_ratings": [
1624
          action: 0,
1625
          rating: 0.1},
1626
1627
          action: 1,
1628
          rating: -0.3},
1629
1630
          action: 2,
          rating: 0.9},
           ... (one entry for each legal move)
1633
1634
          reason: Your detailed reasoning for the chosen action,
          action: 2
1635
1637
          IMPORTANT FORMATTING RULES:
1638
          - Rate each legal move from -1 (terrible) to 1 (excellent)
1639
          - Include ALL legal moves in the move_ratings array
1640
          - The action field should be the index of your chosen move
1641
          - Use valid JSON with proper quotes around all strings
1642
1643
          To win, you need to play the cards in the correct sequence and maximize the total
1644
          score of the fireworks. Good luck!
1645
           Calculate the probability of each card in your hand and the other players' hands to
1646
          make better decisions.
1647
1648
           Card Distribution and Probability Calculation:
1649
          - Each color has a specific number of cards per rank:
1650
          * Rank 1: 3 cards per color (15 total)
1651
          * Rank 2: 2 cards per color (10 total)
1652
          * Rank 3: 2 cards per color (10 total)
          * Rank 4: 2 cards per color (10 total)
          * Rank 5: 1 card per color (5 total)
          - Total deck: 50 cards (5 colors \times 10 cards = 50)
1656
          Try to save the critical cards like rank 5, second card of each color, rank 2,3,4.
1657
1658
           When evaluating unknown cards (your own or others'), calculate probabilities by:
1659
           1. Take the initial distribution of cards and subtract the cards you can see in other players' hands
1660
          2. Subtract cards you can see in the fireworks stacks
1661
          3. Subtract cards that have been discarded (check the discard pile)
1662
          4. Calculate probability
1663
1664
          Use these probability calculations to make better decisions about plays, hints, and dis-
1665
          cards. Make use of the possible cards/ranks provided actively for your decisions and probability
          calculations. They were gathered from historical clues. For example, if you see a card could
          only be green, yellow we can deduce that the card is not red, blue or white. If you see a card
          could only be rank 1, 2, 3 we can deduce that the card is not rank 4 or 5."
1668
1669
           Output:
1671
1672
           "move_ratings": [
```

```
1674
            "action": 0, "rating": 0.7 },
1675
            "action": 1, "rating": 0.7 },
1676
            "action": 2, "rating": 0.7 },
1677
            "action": 3, "rating": 0.8 },
1678
            "action": 4, "rating": 0.7 },
1679
            "action": 5, "rating": -0.5 },
            "action": 6, "rating": -0.5 },
1681
            "action": 7, "rating": -0.5 },
1682
            "action": 8, "rating": -0.6 },
1683
            "action": 9, "rating": -0.5 },
            "action": 10, "rating": -1.0 },
1684
            "action": 11, "rating": -1.0 },
1685
            "action": 12, "rating": -1.0 },
            "action": 13, "rating": -1.0 },
1687
            "action": 14, "rating": -1.0 },
1688
            "action": 15, "rating": -1.0 },
1689
            "action": 16, "rating": -1.0 }
           "reason": "With the deck empty and this as our final turn, there are no remaining playable cards.
          The cards needed to advance the fireworks (G4, G5, W3, W5, B5) are all in the discard pile. Any
1693
          attempt to play a card will result in a penalty. Giving a hint is also pointless, as this is the last
          turn of the game. Therefore, the only logical move is to discard a card to end the game without
1695
          losing a life token. The card in position 3 is known to be Blue, and given the state of the board
          and discards, it is unplayable. Discarding this known useless card is the safest and most optimal
          action.",
           "action": 3
1698
1699
1700
```

We have added additional content from "Critical: Also include move ratings" through to the end of the prompt. If we remove this section, as well as the final round details, the prompt reverts to the SPIN-Bench setup. For the results shown in Figure 2, we further removed the discard pile and the deduction statements respectively (those beginning with phrases like this could be for both the current player and other players).

J BEST OF K - FINAL AGENT'S PROMPT (BOTH MINCON AND DEDUCTON SETUP)

Receives the same input as the single agent setup. Then the following is appended:

Best of K final agent's Prompt

Below are n different responses from the same model to the above game situation. Each response contains reasoning and a chosen move.

```
{Response 1:}
...
{Response n:}
```

1701

1702

1703

1704

1705 1706 1707

1708

1709 1710 1711

17121713

17141715

1716

1717 1718

1719

1722 1723

1724

1725

1726

1727

Our task is to:

- 1. Review all n responses above
- 2. Analyze the reasoning in each response
- 3. Consider which response has the best strategic thinking
- 4. Select the action that you believe is the optimal choice for this game situation

Please provide your reasoning and chosen action in the same format as the responses above.

K MIXTURE-OF-AGENT ARCHITECTURE

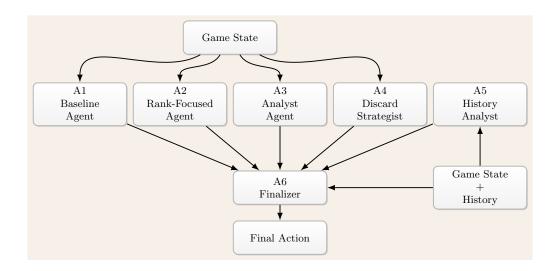


Figure 15: Mixture-of-agent system

L EXAMPLE OF MINCON SETUP MULTI-AGENT PROMPTS:

L.1 SHARED INFORMATION:

This information is common to all agent prompts.

Common Information to all agents

Game State: P0 (5p Game). Lives: 3, Info: 1, Deck: 0.

Fireworks: R4 Y5 G4 W2 B4.

Discards: 1 red card rank 1, 1 red card rank 3, 1 red card rank 4, 1 red card rank 5, 1 yellow card rank 2, 1 yellow card rank 3, 1 green card rank 1, 1 green card rank 2, 1 green card rank 3, 1 green card rank 4, 1 green card rank 5, 1 white card rank 2, 1 white card rank 4, 1 blue card rank 2.

FINAL ROUND! 1 turns left (P0 is last).

Visible Hands:

P1 Hand: [White 5, White 1, Red 2].

P2 Hand: [Yellow 4, White 1, Yellow 1].

P3 Hand: [White 3, Blue 4, White 4, Blue 1].

P4 Hand: [Blue 1, Blue 3, Yellow 1]

Your Knowledge (Hints):

[UnknownColor 3, UnknownColor UnknownRank, UnknownColor UnknownRank, UnknownColor UnknownRank] (Indices 0-3).

Others' Knowledge:

P1 Knows: [UnknownColor UnknownRank, UnknownColor UnknownRank, UnknownColor UnknownRank, [UnknownColor UnknownRank]].

P2 Knows: [UnknownColor 4, UnknownColor UnknownRank, UnknownColor UnknownRank, [UnknownColor UnknownRank]].

P3 Knows: [UnknownColor UnknownRank, UnknownColor UnknownRank, UnknownRank, UnknownColor UnknownRank, Unknown

```
1782
          UnknownRank, UnknownColor UnknownRank].
1783
          P4 Knows: [Blue UnknownRank, Blue UnknownRank, UnknownColor UnknownRank, [Un-
1784
          knownColor UnknownRank]]
1785
          Legal Moves:
1786
          (Discard 0)
1787
          (Discard 1)
1788
          (Discard 2)
1789
          (Discard 3)
1790
          (Play 0)
1791
          (Play 1)
1792
          (Play 2)
          (Play 3)
1793
          (Reveal player +1 color R)
1794
          (Reveal player +1 color W)
1795
          (Reveal player +2 color Y)
1796
          (Reveal player +2 color W)
1797
          (Reveal player +3 color W)
1798
          (Reveal player +3 color B)
1799
          (Reveal player +4 color Y)
1800
          (Reveal player +4 color B)
1801
          (Reveal player +1 rank 1)
1802
          (Reveal player +1 rank 2)
1803
          (Reveal player +1 rank 5)
          (Reveal player +2 rank 1)
          (Reveal player +2 rank 4)
1805
          (Reveal player +3 rank 1)
1806
          (Reveal player +3 rank 3)
1807
          (Reveal player +3 rank 4)
1808
          (Reveal player +4 rank 1)
1809
          (Reveal player +4 rank 3)
1810
          RECENT TURN HISTORY (LAST 10):
1811
1812
                 • T46 (P0, Info:1, FW:R4 Y4 G3 W2 B3): [(Reveal player +2 rank 5)]
1813
                 • T47 (P1, Info:0, FW:R4 Y4 G3 W2 B3): [(Discard 0)]
1814
                 • T48 (P2, Info:1, FW:R4 Y4 G3 W2 B3): [(Reveal player +2 rank 4)]
1815
                 • T49 (P3, Info:0, FW:R4 Y4 G3 W2 B3): [(Discard 0)]
1816
1817
                 • T50 (P4, Info:1, FW:R4 Y4 G3 W2 B3): [(Reveal player +1 rank 4)]
1818
                 • T51 (P0, Info:0, FW:R4 Y4 G3 W2 B3): [(Play 0)]
1819
                 • T52 (P1, Info:0, FW:R4 Y4 G4 W2 B3): [(Discard 0)]
1820
                 • T53 (P2, Info:1, FW:R4 Y4 G4 W2 B3): [(Play 3)]
1821
                 • T54 (P3, Info:2, FW:R4 Y5 G4 W2 B3): [(Reveal player +1 color B)]
1823
                 • T55 (P4, Info:1, FW:R4 Y5 G4 W2 B3): [(Play 3)]
1824
1825
```

AGENT 1 PROMPT:

1826

1827 1828

1829 1830

1831 1832

1833 1834

1835

Everything same as the MinCon single agent setup.

AGENT 2 PROMPT:

Same input as Agent 1 with the following appended to the system prompt:

"with a preference for rank clues over color clues when both are equally valuable."

1836 AGENT 3 (ANALYST) PROMPT: 1837 1838 **System Prompt** You are the Analyst Agent. Your task is to analyze all legal moves and 1840 provide a detailed assessment of their potential value. 1841 YOUR TASK: • For PLAY moves: Assess likelihood of success (Certain, High, Medium, Low, 1843 Impossible). 1844 • For DISCARD moves: Assess safety (High, Medium, Low, Very Low). 1845 For CLUE moves: Evaluate information value (High, Medium, Low). 1847 **OUTPUT FORMAT: Move Analysis:** Move 0 (Type): [Detailed analysis of the move's value and risk] 1849 Move 1 (Type): [Detailed analysis of the move's value and risk] ... (continue for all moves) 1850 **Summary:** 1851 Brief summary of the most promising moves and any key observations", User Prompt You are the Analyst Agent. Analyze all legal moves and provide a detailed assessment of their potential value. [Game State] 1855 [Legal moves] 1857 1860 AGENT 4 (DISCARD STRATEGIST) PROMPT: 1861 1862 1863 System Prompt You are the Discard Pile Analyst. Your task is to analyze the discard pile 1864 and provide insights about what cards are safe to discard based on what has already been 1865 discarded. 1866 YOUR TASK: 1867 1. Discard Pile Analysis: 1868 * Analyze what cards of each color and rank have been discarded * Identify which cards are now impossible to complete their fireworks * Note which high-value cards (5s) or critical cards are already discarded 2. Safe Discard Recommendations: * Based on the discard pile, identify which types of cards would be safe to discard 1872 * Highlight any cards that should absolutely not be discarded due to what's already in the 1873 discard pile 1874 **OUTPUT FORMAT:** 1875 **Discard Pile Status:** 1876 Detailed analysis of what's in the discard pile by color and rank 1877 **Critical Cards Lost:** 1878 List of important cards that are already discarded 1879 Safe Discard Recommendations: 1880 List of card types that would be safe to discard based on the discard pile analysis 1881 **User Prompt** You are the Discard Pile Analyst. Analyze the discard pile and provide insights about what cards are safe to discard. [Game State]

AGENT 5 (HISTORY ANALYST) PROMPT:

[Legal moves]

"system_prompt": "You are Agent 5, a History Analyst. Your task is to analyse the recent turn history in the context of the current game state. Provide concise insights and potential inferences. The user prompt will contain the current Game State and Recent Turn History. FOCUS ON: * Patterns and trends in players' decisions * Inferences about unknown cards based on past plays/clues * Strategic opportunities based on history * Potential warnings or red flags **OUTPUT FORMAT: History Insights:** List of key insights from history that could inform the current decision", "user_prompt": "You are Agent 5 (History Analyst). Analyze the recent turn history in the context of the current game state. Provide concise insights and potential inferences. Do NOT propose a move. [Game State] [Legal moves] [Recent Turn History] Your Task: * Identify any notable patterns (e.g., repeated clues, specific discards). * Infer potential player intentions or card knowledge based on actions. * Highlight any warnings or opportunities suggested by the history. * Keep insights brief and relevant to the *current* decision. **Output Format: History Insights:** - [Insight 1] - [Insight 2] - ..."

AGENT 6 PROMPT:

"system_prompt": "You are Agent 6, the Finalizer Agent in a cooperative Hanabi game. You make the FINAL DECISION based on all other agents' inputs.

The user prompt will contain the Game State, Legal Moves, proposals from other agents,

analysis, and history insights.

1944

1945

1946

1947

1948 1949

1950

1951

1952

1953

1957

1959

1961

1963

1964

1965

1966

1967

1968

1969

1970

1972

1974

1975

1976

1981

1982

1984

1985

1987

1988

1992

1993

1997

Hanabi Strategic Considerations:

- * Playing Cards: Consider playing a card if it's KNOWN (both color and rank) and is the *exact next card needed* for a firework. Such plays are generally very strong. Explain the basis for this knowledge.
- * **Giving Clues:** When information tokens are available (especially if the count is healthy, e.g., > 1-2, unless a clue is critical):
- * Think about clues that could enable another player to make a safe play soon.
- * Consider clues to help save important cards (like unique 5s or cards needed to complete a suit if other copies are gone).
- * Aim for clues that offer new, non-redundant information. Touching multiple cards efficiently can be good. (Always check 'Others' Knowledge' to avoid giving information already known).
- * Assess if the current token count supports giving a clue, especially if it doesn't lead to an immediate play.
- * If a clue seems valuable (high impact, not redundant, affordable), explain its benefits. Otherwise, discarding might be a better option.
- * **Discarding Cards:** If there isn't a clear safe play and giving a valuable clue isn't feasible (or info tokens are at 0):
- * Consider discarding the "safest" card. This could be one known to be useless (e.g., a duplicate of an already played/discarded card, or a card for a completed firework).
- * If no card is known to be useless, think about discarding one with the least information or one deemed least likely to be critical.
- * Explain why the chosen discard is considered the safest. Discarding helps regain information tokens.
- * Do not take unnecesary risk especially if the life token is 1.

DECISION PROCESS:

Your decision should be guided by the Hanabi Strategic Considerations, taking into account all provided inputs. Carefully weigh the options:

- * **Playing a card:** Especially if it's known to be safe and needed.
- * Giving a clue: If it's valuable (enables a play, saves a card, non-redundant) and tokens are sufficient.
- * **Discarding a card:** If playing or cluing isn't a better option, or tokens are critically low. **WEIGH ALL INPUTS:**
- Agent 1 General move suggestions
- Agent 2 Alternative move suggestions
- Agent 3 Detailed hand and clue analysis
- Agent 4 Discard expertise and justification for/against discarding
- Agent 5 History insights, patterns, and inferences

Consider the specific advice from Agent 3 on playability/discard safety and Agent 4's discard recommendation. Agent 5's insights might reveal hidden opportunities or risks.

Evaluate if any card is a known safe play (e.g., Agent 3 indicates Certain playability, or it's self-evident from your knowledge). Such plays are often strong.

If not, carefully compare the potential benefits of the best available clue (considering value assessed by Agent 3 and strategic fit) against the necessity and safety of a discard (considering Agent 3's safety assessment and Agent 4's proposal).

Be cautious with life tokens; risky plays are generally for late-game high potential gain if lives are > 1. Do not give redundant clues. Discarding early can be appropriate if tokens are needed and no clearly better option exists. Protect 5s.

1998 **OUTPUT FORMAT:** 1999 Reasoning: [Your final reasoning, explaining why you chose this move based on the 2000 agents' input and the strategic considerations. Reference specific agent inputs if they were 2001 influential.] 2002 Move Ratings: [Rate EACH legal move from -1 (bad) to 1 (excellent), e.g., Move 0: 0.9, 2003 Move 1: -0.5, Move 2: 0.2, ...] 2004 Chosen Move Number: [number of the best move] 2005 Do not add * before or after Chosen Move Number", 2006 2007 "user_prompt": "You are Agent 6, the Finalizer Agent. Decide the single best move for 2008 the current player. 2009 First, check for KNOWN SAFE PLAYS according to your strict system prompt definition. If one exists, you MUST choose it. If no safe play exists, review the proposals (Agents 1, 2), discard proposal (Agent 4), analyst 2011 assessment (Agent 3: hand & clues), history analysis (Agent 5), and turn history to choose 2012 the best clue or discard. Explain your final reasoning clearly. 2013 2014 [Game State] 2015 [Legal moves] 2016 [Recent Turn History] 2017 2018 Agent 1 Proposal — 2019 [Response A1] 2020 End Agent 1 Proposal — – Agent 2 Proposal – 2021 [Response A2] 2022 — End Agent 2 Proposal — 2023 Agent 3 Analysis (Hand & Clues) — 2024 [Response A3] 2025 — End Agent 3 Analysis — 2026 Agent 4 Discard Proposal — 2027 [Response A4] 2028 — End Agent 4 Discard Proposal — 2029 — Agent 5 History Analysis — 2030 [Response A5] 2031 — End Agent 5 History Analysis — 2032

M Example of DeductCon Setup Multi-agent Prompts:

AGENT 1 PROMPT:

20352036

20372038

2039 2040

20412042

2043 2044

2045

204620472048

2049 2050

2051

Same input as single agent DeductCon prompt setup

AGENT 2 PROMPT:

Same as agent 1 with the following appended to the prompt:

IMPORTANT RULE:

When a color clue and a rank clue are equally valuable, you must give the rank clue.

AN EXAMPLE OF COMMON CONTEXT FOR AGENTS 3, 4, 5 AND 6

This block of text, containing the game rules and the complete, dynamic game state, is prefixed to the instructions for each of the specialist agents.

Common Information

You are a master of hanabi game. You are playing a game of Hanabi with 5 players. Hanabi is a cooperative card game where players work together to create a series of fireworks by playing cards in ascending numerical order starting from 1. Each player holds their cards facing outward so that all players can see everyone else's cards but not their own. The objective is to play cards in sequence (1 through 5) for each color without making mistakes. There are 5 different colors and each color has cards numbered 1 to 5.

Key Rules:

On your turn, you have three types of possible actions:

Give a Hint(Reveal): Provide a hint to another player about their cards, specifying either a color or a number present in their hand. Hints must be accurate and can only reveal positions of cards matching the hint.

Discard a Card: Discard one of your own cards to potentially gain an Info token.

Play a Card: Attempt to play a card from your hand. If played correctly in sequence, it adds to the fireworks; if not, it reduces one life token.

Tokens:

Life Tokens: Deducted when a wrong card is played.

Info Tokens: Used to give clues.

Illegal Moves: Playing a card that cannot be placed properly costs a life token. If life tokens reach zero, the game ends in failure.

Game End: The game ends when all fireworks are completed (perfect score of 25), or when the deck is exhausted and each player has taken one final turn, or when the players run out of life tokens.

State Representation: The game state is represented with the following details:

Life tokens: Number of remaining life tokens.

Info tokens: Number of available information tokens.

Fireworks: Current progress on each firework color (e.g., R1, Y0, G1, W0, B0).

Discards: Cards that have been discarded.

Your Role:

You are one of the players, cooperating with others to maximize the total score of the fireworks (the number of cards correctly played in sequence).

Although you cannot see your own cards, you can see the cards in the hands of your teammates. Use hints, discards, and plays strategically to guide the team towards successful sequences.

Remember, communication is limited to hints about colors or numbers only, and sharing illegal or extraneous information is not allowed. Work together, follow the rules, and aim for the highest cooperative score possible!

Current Game State:

There are 3 life tokens and 0 info tokens remaining.

The fireworks progress: R stack is at 2, Y stack is at 5, G stack is at 3, W stack is at 2, B stack is at 3.

Your hand contains the following cards:

Card 0:

- Known info: 'X1'. Known: rank is 1.
- Could be any of these colors: Red, Yellow, Blue with ranks: 1.

Card 1

- Known info: 'XX'. No hints about this card's color or rank have been given yet.
- Could be any of these colors: Red, Yellow, Green, Blue with ranks: 1, 3.

2106 Card 2: 2107 - Known info: 'X4'. Known: rank is 4. 2108 - Could be any of these colors: Red, Yellow, Green, Blue with ranks: 4. 2109 Card 3: 2110 - Known info: 'XX'. No hints about this card's color or rank have been given yet. 2111 - Could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 5. 2112 From your perspective, you can see the other players' hands clearly. Here's what you 2113 observe: 2114 Player +4's hand: 2115 - A card: You can see the card: 'W4', This player has no specific hints about the card's 2116 identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 4, 5. 2117 - A card: You can see the card: 'Y1', This player has no specific hints about the card's 2118 identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue 2119 with ranks: 1, 2, 4, 5. 2120 - A card: You can see the card: 'R4', This player has no specific hints about the card's 2121 identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue 2122 with ranks: 1, 2, 3, 4, 5. 2123 - A card: You can see the card: 'B4', This player has no specific hints about the card's 2124 identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue 2125 with ranks: 1, 2, 3, 4, 5. 2126 Player +1's hand: 2127 - A card: You can see the card: 'G5', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Green, White, Blue with ranks: 1, 2, 2128 3, 4, 5. 2129 - A card: You can see the card: 'Y2', This player has no specific hints about the card's 2130 identity, This player knows it could be any of these colors: Yellow, Green, White, Blue with 2131 ranks: 1, 2, 3, 4, 5. 2132 - A card: You can see the card: 'R1', This player has no specific hints about the card's 2133 identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue 2134 with ranks: 1, 2, 3, 4, 5. 2135 - A card: You can see the card: 'R2', This player has no specific hints about the card's 2136 identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue 2137 with ranks: 1, 2, 3, 4, 5. Player +2's hand: - A card: You can see the card: 'R5', This player has no specific hints about the card's 2139 identity, This player knows it could be any of these colors: Red, Yellow, Green, Blue with ranks: 2140 2141 2142 2143 3, 4, 5. 2144 2145 2146

- A card: You can see the card: 'G4', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Red, Yellow, Green, Blue with ranks:
- A card: You can see the card: 'Y4', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Red, Yellow, Green, White, Blue with ranks: 1, 2, 3, 4, 5.

Player +3's hand:

- A card: You can see the card: 'W3', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Red, Yellow, White with ranks: 1, 2,
- A card: You can see the card: 'W2', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Red, Yellow, White with ranks: 1, 2,
- A card: You can see the card: 'Y3', This player has no specific hints about the card's identity, This player knows it could be any of these colors: Red, Yellow, White, Blue with ranks: 1, 2, 3, 4, 5.
- There are 0 cards remaining in the deck. The discard pile contains: 2 red cards rank 3, 1 red card rank 4, 2 green cards rank 1, 1 green card rank 2, 1 green card rank 3, 1 green card rank 4, 2 white cards rank 1, 1 white card rank 3, 1 white card rank 4, 1 white card rank 5, 1 blue card

2158 2159

2147

2148

2149

2150

2151

2152 2153

2154

2155

2156

rank 1, 1 blue card rank 2, 1 blue card rank 3, 1 blue card rank 5.

FINAL ROUND: The deck is empty. You are the final player and this is the final turn for the whole game.

```
2214
        AGENT 3 (ANALYST) PROMPT:
2215
2216
2217
2218
            [Shared Context]
2219
           Analyse EVERY candidate move based on the game state provided above.
2220
2221
           Legal Moves:
2222
            "0": "(Discard 0)",
2223
           "1": "(Discard 1)",
2224
           "2": "(Discard 2)",
2225
           "3": "(Discard 3)",
2226
           "4": "(Play 0)",
2227
           "5": "(Play 1)",
2228
           "6": "(Play 2)",
2229
           "7": "(Play 3)"
2230
2231
2232
           For EVERY move listed above, provide a structured analysis using the following
2233
           template. Be detailed.
2234
           Move 0:
2235
           Type: <Play / Discard / Color-Clue / Rank-Clue>
2236
           Reason: ...
2237
           Immediate_effect: <score change, token gain/loss, or no immediate change>
2238
           Reason: ...
2239
           Probability_of_success: <Certain / High / Medium / Low / Impossible>; for plays
2240
           Reason: . .
2241
           Discard_risk_level: <Very-Safe / Safe / Risky / Deadly>; for discards
2242
           Reason: ...
2243
           Clue_value: <Immediate-Play / Critical-Save / Setup / Redundant / Wasted>; for clues
2244
           Reason: ...
           Info token cost or gain: <+1/0/-1>
2245
           Reason: ...
2246
           Future_impact: <detailed sentence on longer-term effect.>
2247
           Overall_rationale: <integrate all factors above.>
2248
2249
           (repeat this full block for EVERY legal move)
2250
2251
           Summary:
2252
           2253
           Major_risks_detailed: <paragraph on biggest dangers.>
2254
           2255
           Calculate the probability of each card in your hand and the other players' hands to
2256
           make better decisions.
2257
2258
           Card Distribution and Probability Calculation
2259
           - Each color has a specific number of cards per rank:
2260
           * Rank 1: 3 cards per color (15 total)
2261
           * Rank 2: 2 cards per color (10 total)
2262
           * Rank 3: 2 cards per color (10 total)
2263
           * Rank 4: 2 cards per color (10 total)
2264
           * Rank 5: 1 card per color (5 total)
2265
           - Total deck: 50 \text{ cards } (5 \text{ colors } \times 10 \text{ cards} = 50)
2266
2267
```

Try to save the critical cards like rank 5, second card of each color, rank 2,3,4.

When evaluating unknown cards (your own or others'), calculate probabilities by:

- 1. Take the initial distribution of cards and subtract the cards you can see in other players' hands
- 2. Subtract cards you can see in the fireworks stacks
- 3. Subtract cards that have been discarded (check the discard pile)
- 4. Calculate probability

Use these probability calculations to make better decisions about plays, hints, and discards. Make use of the possible cards/ranks provided actively for your decisions and probability calculations. They were gathered from historical clues. For example, if you see a card could only be green, yellow we can deduce that the card is not red, blue or white. If you see a card could only be rank 1, 2, 3 we can deduce that the card is not rank 4 or 5.

2322 AGENT 4 (DISCARD STRATEGIST) PROMPT: 2324 2325 2326 [Shared Context] 2327 For EVERY card in the current player's hand, provide a detailed discard analysis based on 2328 the game state above. 2330 Card 0: 2331 Safety_probability: <0-1> 2332 Reason: ... 2333 Criticality: <Very-High / High / Medium / Low / Very-Low> 2334 Reason: ... Visible_duplicates: "X of Y copies seen – location(s): ..." (If there are no visible duplicates, 2335 write "None") 2336 Reason: ... 2337 Recommendation: <Discard / Keep> 2338 Reason: ... 2339 2340 (repeat for all cards in the hand) 2341 2342 Detailed_Summary: 2343 Safest discards: <paragraph naming the safest card(s) and why.> Cards_to_protect: <paragraph naming risky cards and why.> 2345 Distribution_notes: <paragraph noting colours/ranks exhausted or at single copy.> 2346 Like firework red is already at 3, Two red 4 is already in the discard pile so we can 2347 discard the red card in our hand. 2348 2349 2350 Calculate the probability of each card in your hand and the other players' hands to 2351 make better decisions. 2352 2353 Card Distribution and Probability Calculation 2354 - Each color has a specific number of cards per rank: 2355 * Rank 1: 3 cards per color (15 total) * Rank 2: 2 cards per color (10 total) 2356 * Rank 3: 2 cards per color (10 total) 2357 * Rank 4: 2 cards per color (10 total) * Rank 5: 1 card per color (5 total) 2359 - Total deck: $50 \text{ cards } (5 \text{ colors } \times 10 \text{ cards} = 50)$ 2360 2361 Try to save the critical cards like rank 5, second card of each color, rank 2,3,4. 2362 2363 When evaluating unknown cards (your own or others'), calculate probabilities by: 2364 1. Take the initial distribution of cards and subtract the cards you can see in other players' 2365 2366 2. Subtract cards you can see in the fireworks stacks 3. Subtract cards that have been discarded (check the discard pile) 2367 4. Calculate probability 2368 2369 Use these probability calculations to make better decisions about plays, hints, and 2370 discards. Make use of the possible cards/ranks provided actively for your decisions and 2371 probability calculations. They were gathered from historical clues. For example, if you see a 2372 card could only be green, yellow we can deduce that the card is not red, blue or white. If you 2373 see a card could only be rank 1, 2, 3 we can deduce that the card is not rank 4 or 5. Use this 2374 to Backup your decision to discard or save a card.

AGENT 5 (HISTORY ANALYST) PROMPT:

Your identity for this turn is Player 1 (P1).

[Shared context]

IMPORTANT: In the history below, when you see a clue like '(Reveal player +2 color R)', the '+2' refers to the position relative to the player who GAVE the clue, not relative to you (the current player). For example, if Player +1 gave a clue to Player +3, it means they clued the player who is 2 positions ahead of them.

Turn 1: Player +2 (P3) chose move '(Reveal player +4 rank 1)'. Fireworks: R0, Y0, G0, W0, B0 \rightarrow R0, Y0, G0, W0, B0, Info tokens: $8\rightarrow$ 7.

Turn 2: Player +3 (P4) chose move '(Reveal player +1 rank 1)'. Fireworks: R0, Y0, G0, W0, $B0\rightarrow R0$, Y0, G0, W0, B0, Info tokens: $7\rightarrow 6$.

... (full history from Turn 3 to 57) ...

Turn 58: Player +4 (P0) chose move '(Reveal player +1 rank 4)'. Fireworks: R2, Y5, G3, W2, B3 \rightarrow R2, Y5, G3, W2, B3, Info tokens: $1\rightarrow$ 0.

For relevant turns above, explain what the acting player was trying to achieve and what that reveals about hidden cards. (Mostly focus on recent turns and think why would someone give clues to other players instead of giving clue to us? or why someone prioritise us over other players? The same with different cards in our hand.)

Speculations:

- player+4 gave me a Yellow-colour clue instead of clueing player+1's Yellow card while the Yellow stack is at 3. Yellow 1 and Yellow 3 are already in the discard pile, so my hidden card can only be Yellow 2 or Yellow 4. Because a Yellow 2 would not score immediately, the clue strongly implies my card is Yellow 4 and ready to play.
- player+1 did not clue my right-most card even though it could be playable next if it were Red 2. That suggests they believe it is not Red 2, increasing the likelihood that my left-most card (just clued) is the immediate scoring card.

Calculate the probability of each card in your hand and the other players' hands to make better decisions.

Card Distribution and Probability Calculation

- Each color has a specific number of cards per rank:
- * Rank 1: 3 cards per color (15 total)
- * Rank 2: 2 cards per color (10 total)
- * Rank 3: 2 cards per color (10 total)
- * Rank 4: 2 cards per color (10 total)
- * Rank 5: 1 card per color (5 total)
- Total deck: 50 cards (5 colors x 10 cards = 50)

Try to save the critical cards like rank 5, second card of each color, rank 2,3,4.

When evaluating unknown cards (your own or others'), calculate probabilities by:

- 1. Take the initial distribution of cards and subtract the cards you can see in other players' hands
- 2. Subtract cards you can see in the fireworks stacks
- 3. Subtract cards that have been discarded (check the discard pile)
- 4. Calculate probability

Use these probability calculations to make better decisions about plays, hints, and discards. Make use of the possible cards/ranks provided actively for your decisions and probability calculations. They were gathered from historical clues. For example, if you see a card could only be green, yellow we can deduce that the card is not red, blue or white. If you see a card could only be rank 1, 2, 3 we can deduce that the card is not rank 4 or 5. Use this to backup your speculations.

AGENT 6 PROMPT:

```
2432
            [Shared Context]
2433
2434
2435
           You have also received:
2436
           - Ratings JSON from the first strategist
2437

    Ratings JSON from the rank-preferring strategist

2438
           - Full move analysis text
2439

    Discard-probability report

2440
           - History deductions text
2441
            Recent Game History:
2442
            [Recent Game History]
2443
2444
           Report from Agent 1 (Baseline):
2445
            [Response from A1]
2446
2447
           Report from Agent 2 (Rank-Preferring):
2448
            [Response from A2]
2449
2450
            Report from Agent 3 (Analyst):
2451
            [Response from A3]
2452
           Report from Agent 4 (Discard Expert):
2453
            [Response from A4]
2454
2455
           Report from Agent 5 (Historian):
2456
            [Response from A5]
2457
2458
2459
            Combine all of that and choose the single best move. Your output must be a single,
2460
           valid JSON object.
2461
2462
            "move_ratings": [ ... include every legal move with a rating –1 to 1 ... ],
2463
            "reason": "short justification that cites insights from earlier analyses",
2464
            "action": <index of chosen move>
2465
2466
2467
2468
```

N MULTI-TURN

N.1 Multi-turn Prompts

```
2484
             "action": 6, "rating": -0.8},
2485
            {"action": 7, "rating": -0.8}
2486
2487
            "deduction": {
2488
            "player+1": {
            "card0": "color could be Y,G,W,B; rank could be 1,2,4,5",
2489
            "card1": "color could be R,Y,G,W,B; rank could be 1,2,4,5",
2490
            "card2": "color could be R,Y,G,W,B; rank could be 1,2,4,5"
2491
            "card3": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2492
2493
            "player+2": {
            "card0": "color could be R,Y,G,W; rank could be 1,2,3,4,5",
2494
            "card1": "color could be R,Y,G,W; rank could be 1,2,3,4,5",
2495
            "card2": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5",
2496
            "card3": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2497
2498
             "player+3": {
            "card0": "color could be R,Y,G,W; rank could be 1,3,5",
2499
            "card1": "color could be R,Y,G,W; rank could be 1,2,3,5"
            "card2": "color could be R,Y,G,W,B; rank could be 1,2,3,5"
2501
            "card3": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2503
            "player+4": {
            "card0": "color could be R,Y,G,W,B; rank could be 1,3,4,5",
            "card1": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2505
            "card2": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2506
            "card3": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2507
2509
            "reason": With no info tokens available, I cannot give hints and must choose between discarding or playing
            a card. I deduced that my card 0 is rank 3 and must be B3, as all other rank 3 cards are accounted for in the
2510
            played stacks, discarded pile, or other players' hands. Since the blue firework is at 2, B3 is playable and
2511
            safe to play, scoring a point and advancing the firework. Playing card 0 is low-risk and beneficial, especially
2512
            with only one card left in the deck and limited turns remaining. Discarding would gain an info token but is
2513
            suboptimal as it delays scoring and the token may not be used effectively by others before the game ends.
2514
            Other plays are risky due to uncertain playability, and discarding known useful cards is detrimental.
            "action": 4
2515
2516
            You are a master of hanabi game. You are playing a game of Hanabi with 5 players. Hanabi is a cooperative
2517
            card game where players work together to create a series of fireworks by playing cards in ascending
2518
            numerical order starting from 1. Each player holds their cards facing outward so that all players can see
2519
            everyone else's cards but not their own. The objective is to play cards in sequence (1 through 5) for each
            color without making mistakes. There are 5 different colors and each color has cards numbered 1 to 5.
2521
            Key Rules:
2522
2524
```

On your turn, you have three types of possible actions:

2538

2539 2540

Give a Hint(Reveal): Provide a hint to another player about their cards, specifying either a 2541 color or a number present in their hand. Hints must be accurate and can only reveal positions 2542 of cards matching the hint. 2543 Discard a Card: Discard one of your own cards to potentially gain an Info token. 2544 Play a Card: Attempt to play a card from your hand. If played correctly in sequence, it adds 2545 to the fireworks; if not, it reduces one life token. 2546 2547 Tokens: 2548 Life Tokens: Deducted when a wrong card is played. 2549 Info Tokens: Used to give clues. Illegal Moves: Playing a card that cannot be placed properly costs a life token. If life tokens 2550 reach zero, the game ends in failure. 2551 Game End: The game ends when all fireworks are completed (perfect score of 25), or when 2552 the deck is exhausted and each player has taken one final turn, or when the players run out of 2553 life tokens. 2554 2555 State Representation: The game state is represented with the following details: 2557 Life tokens: Number of remaining life tokens. Info tokens: Number of available information tokens. 2559 Fireworks: Current progress on each firework color (e.g., R1, Y0, G1, W0, B0). 2560 Discards: Cards that have been discarded. 2561 Your Role: 2563 You are one of the players, cooperating with others to maximize the total score of 2564 the fireworks (the number of cards correctly played in sequence). 2565 Although you cannot see your own cards, you can see the cards in the hands of your 2566 teammates. 2567 Use hints, discards, and plays strategically to guide the team towards successful sequences. 2568 2569 Remember, communication is limited to hints about colors or numbers only, and 2570 sharing illegal or extraneous information is not allowed. Work together, follow the rules, and 2571 aim for the highest cooperative score possible! 2572 Please think step by step based on the current state 2573 2574 # Think step by step 2575 2576 ## Evaluate Playable Cards in Hand 2577 2578 Look at each card in your hand. 2579 Cross-reference with the current game state to see if any card can be immediately played to 2580 complete or extend a firework stack. 2581 Consider hints you have received about each card (color/rank information) to determine if it 2582 might be safe to play. If a card can be played without risk, prioritize playing it to score a point. 2583 2584 ## Consider Teammates' Hands and Hint Opportunities 2585 2586 Analyze the visible cards in your teammates' hands. 2587 Identify if any of their cards can now be played based on the current firework stacks or 2588 previous hints. 2589 If you notice a teammate holds a card that can be played but they may not realize it, think about what hints you could give them.

Use hints to communicate critical information, such as color or rank, to help them make the right play.

Choose the hint that maximizes the chance for a correct play while considering the limited

Assess Discard Options to Gain Info Tokens

Look for cards in your hand that are least likely to be playable or helpful in the near future.

Consider the remaining deck composition and cards already played/discarded to predict the value of each card.

Discard a card that you believe to be least useful to gain an Info token, especially if no immediate playable or hint options are available.

Ensure that discarding this card won't permanently remove a critical card needed to complete any firework stack.

Now it's your turn. You can choose from the following legal actions:

The legal actions are provided in a mapping of action identifiers to their descriptions:

```
{legal_moves_dict}
```

(Reveal player +N color C): Give a hint about color C to the player who is N positions ahead of you.

(Reveal player +N rank R): Give a hint about rank R to the player who is N positions ahead. (Play X): Play the card in position X from your hand (Card 0, Card 1, Card 2, etc.).

(Discard X): Discard the card in position X from your hand (Card 0, Card 1, Card 2, etc.).

Based on the annotated state and the list of legal actions, decide on the most appropriate move to make. Consider factors like current tokens, firework progress, and information available in hands. Then, output one of the legal action descriptions as your chosen action.

Your output should be in this format:

{"reason": string, "action": int} And the action should be one of the legal actions provided above.

You can only use json valid characters. When you write json, all the elements (including all the keys and values) should be enclosed in double quotes!!!

CRITICAL: Also include move ratings and deduction of what others know about their cards based on the history in this exact JSON format:

```
{
"move_ratings": [
{"action": 0, "rating": 0.1},
{"action": 1, "rating": -0.3},
{"action": 2, "rating": 0.9},
... (one entry for each legal move)
],
"deduction":["player+1": {card1: color is .. or color cannot be . rank is .. or rank cannot be. card2: ....},
"player+2": {....} and so on ]
"reason": Your detailed reasoning for the chosen action;
"action": 2
}
```

CRITICAL: You MUST generate the 'deduction' block by meticulously tracking what each player knows about their own hand. Follow this exact, step-by-step logic for EVERY player on EVERY turn:

2646 Definition: The 'deduction' field must track the accumulated knowledge a player has about 2647 their own cards by listing all remaining possibilities for 'color' and 'rank'. This is built from 2648 the complete public history of hints and actions. 2649 2650 Deduction Logic (Follow these steps for each player): 2651 2652 1. Recall Previous State: Start with the list of possibilities for each card from the previous 2653 turn. (For Turn 1, all cards start with "color could be R, Y, G, W, B; rank could be 1, 2, 3, 4, 5"). 2654 2655 2. Analyze the Most Recent Action: Look at the last move made before your turn. 2656 2657 * If a Hint was GIVEN TO this Player: * Update with Positive Information: For the card(s) identified by the hint, narrow down the 2658 possibilities. If the hint was Blue, the deduction for that card's color becomes "color is Blue." 2659 * Update with Negative Information (MANDATORY): For all other cards in their hand not 2660 identified by the hint, you MUST remove the hinted value from their list of possibilities. 2661 (e.g., color possibilities become R, Y, G, W). 2663 * If this Player ACTED (Played or Discarded): * This is a critical state update. Follow this sequence carefully: 2665 * The card they acted on is removed from their hand. * Retain Knowledge: For all other cards remaining in their hand, their known information is 2667 retained, but their position shifts to fill the gap. 2668 * The new card drawn into the last slot of their hand is a complete unknown. Its deduction is: "color could be R, Y, G, W, B; rank could be 1, 2, 3, 4, 5." 2669 2670 3. Synthesize and Format: Present the final list of possibilities for each card in its 2671 new position. 2672 2673 Example of Correct Deduction: 2674 2675 * Scenario: Player+1 has a hand of R2, B4, W2. It is your turn. In the previous 2676 round, another player gave Player+1 a rank 2hint. 2677 * Your Deduction Output for Player+1 MUST be: 2678 ```json 2679 "player+1": { 2680 "card0": "color could be R, Y, G, W, B; rank is 2; 2681 "card1": "color could be R, Y, G, W, B; rank could be 1, 3, 4, 5", "card2": "color could be R, Y, G, W, B; rank is 2" 2683 2684 2686 Example of a Player Action (Play/Discard): 2687 2688 * Scenario: It is Turn 5. On Turn 4, Player+1 had the following knowledge about 2689 their 4-card hand: 2690 * card0: "color could be R, Y, G, W, B; rank is 2" * card1: "color is Blue; rank could be 3, 4" 2691 * card2: "color could be R, Y, G, W, B; rank is 5" 2692 * card3: "color could be Y, G, W, B; rank could be 1, 2, 3, 4, 5" (They were previously told 2693 their other cards were not Red) 2694 2695 * Action: On their turn, Player+1 plays card 1.

* Your Deduction Output for Player+1 on Turn 5 MUST be:

2697

```
2700
                `ison
2701
             "player+1": {
2702
             "card0": "color could be R, Y, G, W, B; rank is 2,
2703
             "card1": "color could be R, Y, G, W, B; rank is 5",
2704
2705
             "card2": "color could be Y, G, W, B; rank could be 1, 2, 3, 4, 5",
2706
             "card3": "color could be R, Y, G, W, B; rank could be 1, 2, 3, 4, 5"
2708
2709
2710
```

(Notice how the knowledge for the old card 0 remains at position 0, the knowledge for the old card 2 shifts to position 1, the knowledge for the old card 3 shifts to position 2, and the new card at position 3 is completely unknown).

Do not be lazy. You MUST perform this full analysis for all four other players and all of their cards to ensure the 'deduction' block is 100% accurate. An incorrect deduction state will lead to poor team performance.

IMPORTANT FORMATTING RULES:

- Rate each legal move from -1 (terrible) to 1 (excellent)
- Include ALL legal moves in the move_ratings array
- The "action" field should be the index of your chosen move
- Use valid JSON with proper quotes around all strings

To win, you need to play the cards in the correct sequence and maximize the total score of the fireworks. Good luck!

Calculate the probability of each card in your hand and the other players' hands to make better decisions.

Card Distribution and Probability Calculation:

- Each color has a specific number of cards per rank:
- * Rank 1: 3 cards per color (15 total)
- * Rank 2: 2 cards per color (10 total)
 - * Rank 3: 2 cards per color (10 total)
 - * Rank 4: 2 cards per color (10 total)
 - * Rank 5: 1 card per color (5 total)
 - Total deck: $50 \text{ cards } (5 \text{ colors} \times 10 \text{ cards} = 50)$

Try to save the critical cards like rank 5, second card of each color, rank 2,3,4.

When evaluating unknown cards (your own or others'), calculate probabilities by:

- 1. Take the initial distribution of cards and subtract the cards you can see in other players' hands
- 2. Subtract cards you can see in the fireworks stacks
- 3. Subtract cards that have been discarded (check the discard pile)
- 4. Calculate probability

Use these probability calculations to make better decisions about plays, hints, and discards. Make use of the possible cards/ranks provided actively for your decisions and probability calculations. They were gathered from historical clues. For example, if you see a card could only be green, yellow we can deduce that the card is not red, blue or white. If you see a card could only be rank 1, 2, 3 we can deduce that the card is not rank 4 or 5.

Except for the first turn ever for you, you will receive previous one turn prompt and your reasoning before use that to identify the game sate representation in your previous turn and deduce things using the history happend in the last turn after your played.

Below is the current detailed state information.

```
2754
            Game State:
2755
            You are Player P4, Turn 58
2756
            Since your last turn the following actions occurred:
2757
2758

    P0 (Discard 0) | Fireworks: R3 Y2 G4 W2 B3 | Info: 1

2759
            - P1 (Reveal player P2 color B) | Fireworks: R3 Y2 G4 W2 B3 | Info: 0
2760
            - P2 (Play 2) | Fireworks: R3 Y2 G4 W2 B4 | Info: 0
2761
            - P3 (Play 0) | Fireworks: R3 Y2 G4 W2 B4 | Info: 0
2762
2763
            There are 1 life tokens and 0 info tokens remaining.
2764
            The fireworks progress: R stack is at 3, Y stack is at 2, G stack is at 4, W stack is at 2, B
2765
            stack is at 4.
            Your hand (what you know):
2766
            Card 0: unknown color, rank 4
2767
            Card 1: unknown
2768
            Card 2: unknown
2769
            Card 3: unknown
2770
            From your perspective, you can see the other players' hands clearly. Here's what you
2771
            observe:
2772
            Player +1's hand:
2773
            - G2
2774
            - W4
2775
            - W1
2776
            Player +2's hand:
            - B1
2777
            - W4
2778
            - W2
2779
            - Y4
2780
            Player +3's hand:
2781
            - G3
2782
            - R2
2783
            - R3
2784
            Player +4's hand:
2785
            - B2
            - Y3
2786
            - W1
2787
            There are 0 cards remaining in the deck. The discard pile contains: 2 red cards rank 1, 1 red
2788
            card rank 4, 1 yellow card rank 1, 1 yellow card rank 2, 1 yellow card rank 3, 1 yellow card
2789
            rank 4, 1 yellow card rank 5, 2 green cards rank 1, 1 green card rank 4, 1 green card rank 5, 2
2790
            white cards rank 3, 1 white card rank 5, 1 blue card rank 1, 1 blue card rank 3, 1 blue card
2791
            rank 5.
2792
2793
            FINAL ROUND: The deck is empty. You are the final player and this is the final
2794
            turn for the whole game.
2795
2796
            Please note down all the deductions you make so that they will help you in future
2797
            turns.
2798
            For example, if you have deduced that a card cannot be yellow based on previous
2799
            clues (e.g., when a yellow clue is given to the player and the card is not chosen, you know it
            is not yellow—so only RGBW are possible), or if you know a card cannot be rank 2 or 5
            based on previous clues, write this down.
2802
```

Write down all such deductions in the scratch pad, as they will be useful in future

2804

2805

turns.

```
2808
             Use the history well. Avoid giving the same clue to the same player it would be redundant.
            Try to write down what other players know from the history and your previous turns'
2810
            reasoning, and use this to plan your actions.
2811
2812
            Legal moves this turn:
2813
2814
            0: ((Discard 0))
2815
            1: ((Discard 1))
2816
             2: ((Discard 2))
2817
            3: ((Discard 3))
2818
            4: ((Play 0))
2819
             5: ((Play 1))
            6: ((Play 2))
2820
            7: ((Play 3))
2821
2822
             Output
             "move ratings": [
2825
            {"action": 0, "rating": -0.5},
            {"action": 1, "rating": -0.5},
2827
             "action": 2, "rating": -0.5},
             "action": 3, "rating": -0.5},
2829
              "action": 4, "rating": 0.2},
              "action": 5, "rating": -0.7},
2830
2831
              "action": 6, "rating": -0.7},
             "action": 7, "rating": -0.7}
2832
2833
             "deduction": {
             "player+1": {
             "card0": "color could be R,Y,G,W,B; rank could be 1,2,4,5",
2836
             "card1": "color could be R,Y,G,W,B; rank could be 1,2,4,5".
2837
             "card2": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2838
2839
             "player+2": {
             "card0": "color could be R,Y,G,W; rank could be 1,2,3,4,5",
             "card1": "color could be R,Y,G,W; rank could be 1,2,3,4,5",
2841
             "card2": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2842
             "card3": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2843
             "player+3": {
2845
             "card0": "color could be R,Y,G,W; rank could be 1,3,5",
2846
             "card1": "color could be R,Y,G,W; rank could be 1,2,3,5"
2847
             "card2": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2848
2849
             "player+4": {
2850
             "card0": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5",
2851
             "card1": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5",
2852
             "card2": "color could be R,Y,G,W,B; rank could be 1,2,3,4,5"
2853
             "reason": It is the final turn with no info tokens and an empty deck, so no hints can be given.
2855
             The only way to potentially increase the score from 15 is to play a correct card. My card 0 is
2856
             known to be rank 4, and there is a 20% chance it is R4, which is playable, offering a small
2857
            expected gain in score. Playing other cards is riskier with no information, and discarding gains
            an useless info token and removes a potentially useful card, providing no benefit. Therefore,
2859
            I choose to play card 0 to maximize the slim chance of scoring an additional point.;
             "action": 4
```

N.2 MULTI-TURN RESULTS:

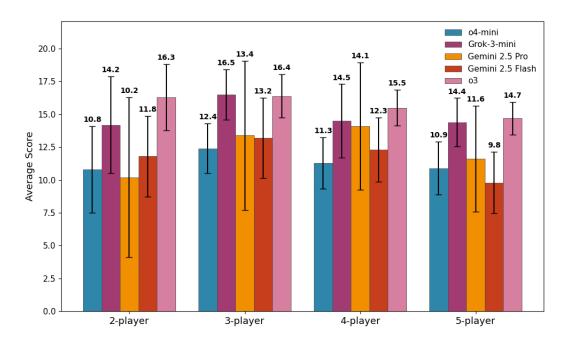


Figure 16: Multi-turn scores of reasoning models across 2-5 player settings.