# Werewolf Arena: A Case Study in LLM Evaluation via Social Deduction

**Suma Bailis,**\* **Jane Friedhoff,** and **Feiyang Chen**

Google Research

`{sbailis,jfriedhoff,feiyangc}@google.com`

## Abstract

This paper introduces Werewolf Arena, a novel framework for evaluating large language models (LLMs) through the lens of the classic social deduction game, *Werewolf*. In Werewolf Arena, LLMs compete against each other, navigating the game's complex dynamics of deception, deduction, and persuasion. The framework introduces a dynamic turn-taking system based on bidding, mirroring real-world discussions where individuals strategically choose when to speak. We demonstrate the framework's utility through an arena-style tournament featuring Gemini and GPT models. Our results reveal distinct strengths and weaknesses in the models' strategic reasoning and communication. These findings highlight Werewolf Arena's potential as a challenging and scalable LLM benchmark.

## 1 Introduction

Creating truly human-like AI requires sophisticated cognitive abilities such as reasoning about others' intentions, navigating deceptive information, and convincingly communicating in complex social settings. Evaluating these nuanced skills in both humans and LLMs poses a significant challenge, as traditional benchmarks often fall short (Ullman, 2023).

Social Deduction Games (SDGs), such as the popular game *Werewolf*, present a compelling avenue for addressing this challenge. SDGs encapsulate many of the complexities of human social interaction, requiring players to utilize various reasoning skills (e.g., temporal, deductive, and inductive) within an uncertain environment. In *Werewolf*, Villager and Werewolf players engage in a battle of wits, leveraging deception and persuasion to achieve their respective goals.

Furthermore, the inherent information asymmetry of *Werewolf*, where only some players possess incomplete knowledge of others' roles, mirrors dynamics of real-world social interactions.

Our Monte Carlo simulation (Algorithm 1) highlights the critical role of strategic communication in overcoming this asymmetry. Without it, Villagers win a mere 1.2% of 100,000 simulated games.

Recently, *Werewolf* has emerged as a popular sandbox for LLM research. Prior work has demonstrated sophisticated agents capable of adapting communication strategies (Jin et al., 2024), employing deductive reasoning for optimal action selection (Xu et al., 2024; Shibata et al., 2023), and using system-2 thinking during gameplay (Wu et al., 2024). Researchers have also explored techniques like retrieval and reflection mechanisms (Xu et al., 2023) to enhance agent learning and created valuable datasets comprising gameplay logs and multimodal artifacts (Lai et al., 2022) to train these agents and establish benchmarks (Chern et al., 2024).

Building on this foundation, Werewolf Arena makes two key contributions.

First, recognizing the importance of strategic communication, we introduce a dynamic turn-taking system where players bid to speak, rather than relying on predefined or random speaking orders. This bidding mechanic closely mirrors real-world group discussions, where individuals strategically time their contributions. This adds a nuanced layer to agent evaluation, allowing us to assess not only *what* is said but also *when* an agent chooses to speak.

Second, unlike past work of developing techniques to improve agent performance, we use *Werewolf* as a proving ground. To evaluate the relative skills of LLMs we have the models play each other. We design a balanced framework where a single model, playing both Villager and Werewolf roles, results in a relatively even win rate for both sides. This balanced setup ensures fair comparisons by minimizing inherent advantages for either

---
\*Corresponding author

**Game settings**

- **Player count:** 8
- **Role distribution:**
  - 🐺 2 Werewolves
  - 🩺 1 Doctor
  - 🔮 1 Seer
  - 4 Villagers
- **Debate length:** 8 turns

**Game over conditions**

**Werewolves win if:**
# Wolves ≥ non-Wolves
**Villagers win if:**
both Werewolves exiled

**(1) 🌙 Night**

**(a) 👀 Special roles execute simultaneously**

| 🐺 Werewolf | 🩺 Doctor | 🔮 Seer |
|---|---|---|
| One Werewolf chooses a player to eliminate. | If still active, protects someone from Werewolves | If still active, investigates someone |

**(b) 📢 Announce outcome**

*If game not over*

**(2) ☀️ Day**

**(c) 💬 Debate**

| 💰 Bid | 🗣️ Speak |
|---|---|
| Players bid to speak next. | Winner speaks. |

**(d) 🗳️ Vote**

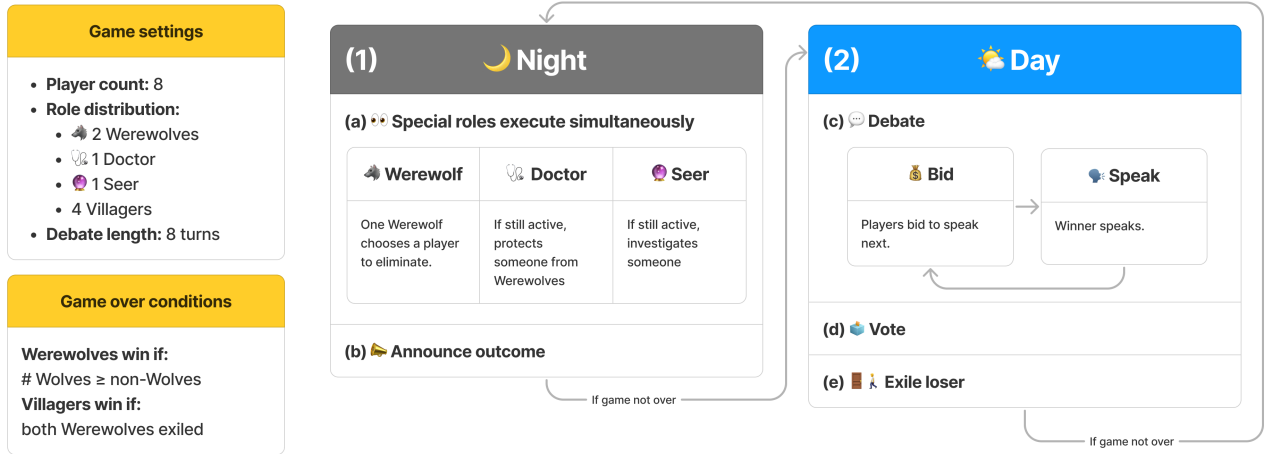**(e) 🚪🚶 Exile loser**

*If game not over*

Figure 1: Game loop of *Werewolf*.

role. Moreover, to facilitate future research and allow others to test their own models, we have released Werewolf Arena at `https://github.com/google/werewolf_arena`.

## 2 Related Work

**Simulating Social Interaction and Strategic Reasoning:** LLMs are transforming agent-based social simulations, enabling agents to communicate, reason (Zhao et al., 2023), solve problems (Wei et al., 2023a), and plan strategically (Song et al., 2023). This has led to advancements in simulating nuanced human behavior in complex social settings (Sreedhar and Chilton, 2024; Zhou et al., 2024; Park et al., 2023; Vezhnevets et al., 2023), and allowed for more believable and capable non-player characters (NPCs) (Ammanabrolu and Riedl, 2019; Urbanek et al., 2019; Wang et al., 2023a). In this landscape, social deduction games, now offer a compelling avenue for studying social dynamics, cooperation, and deception (Kopparapu et al., 2022; Oertel and Salvi, 2013; Vázquez et al., 2015; Leibo et al., 2017; Wang et al., 2023b; Ibraheem et al., 2022).

**Open-Ended Benchmarks for LLMs:** While static reasoning benchmarks are plentiful (et al., 2023; Liang et al., 2023; Hendrycks et al., 2021; Valmeekam et al., 2024), fewer focus on dynamic, competitive evaluations of LLMs. Platforms like LMSYS rely on human evaluation to rank chatbots (Chiang et al., 2024), while Kaggle Simulations, though providing game-like environments, limit games to agents powered by the same model (Kaggle Inc., 2024).

There is room for new benchmarks that evaluate LLMs on their ability to leverage cooperation, deception, and strategic communication in dynamic competition with other models. This type of benchmark offers several benefits: it bypasses the need for human annotations, prevents future data contamination (Deng et al., 2024), and remains relevant as models improve.

## 3 *Werewolf* Environment

This section describes the simulated environment of *Werewolf* we use in Werewolf Arena.

### 3.1 Game Implementation

As illustrated in Figure 1, the game starts with 8 players, consisting of 1 Seer, 1 Doctor, 2 Werewolves, and 4 Villagers. It progresses through rounds until either all Werewolves are exiled (Villager win) or their numbers equal those of the Villagers (Werewolf win). For each game, we randomly select 8 names from a pool of 17 names, to minimize any initial name bias.

Gameplay requires players to discern others' roles while protecting their own identities. The game proceeds in two phases. During each Night, special roles happen simultaneously: the Werewolves conspire to eliminate a single Villager, the Doctor chooses someone to protect, and the Seer investigates a player to learn their role. The Daytime proceeds sequentially: it consists of a structured debate among all players and a subsequent voting session, where a majority is required for exile.

Currently, debates are capped at 8 turns, ensuring each player a chance to speak in the first round.

Future iterations could implement a more dynamic system, ending debates upon consensus (explored in Appendix C). This would address the occasional repetition observed in later rounds with fewer players.

The game is orchestrated via a rules-based Game Master (GM), who oversees the game's progression, assigns roles, and ensures the timely execution of agent actions. The GM also diligently tracks and updates observable game events, like eliminations or exiles, in the agents' memories, keeping them informed of the game's progress.

## 3.2 Agent Architecture

Agents are equipped to perform a suite of actions essential to *Werewolf*'s gameplay:

- **Core Actions**: All agents engage in *voting* to determine player exiles, *debating* to influence others and gather information, and *bidding* for their turn to speak, reflecting the dynamic nature of group discussions.

- **Special Role Actions**: Agents assigned as Werewolves, Doctors, or Seers execute nighttime actions of *eliminating* a villager, *protecting* a player, and *investigating* a player's true role, respectively.

- **Agent Memory**: Drawing inspiration from (Park et al., 2023), each agent possesses a memory stream that contains observational and reflective memories. The observational memories record all game-level events and privileged information accessible to each player based on their role (e.g., a Seer's memory would include the results of their investigations). At the end of each round, agents engage in *summarizing*, distilling key insights from the debate. These reflective summaries enable agents to recall pertinent information and notice patterns in subsequent rounds.

Each action is guided by a tailored prompt template[1]. This template incorporates the agent's memories and the current game state from their perspective, ensuring contextually appropriate actions. To counter the early-game tendencies observed in (Xu et al., 2024) of selecting the first or last option in a list, we randomize the order that player names are presented during voting and special actions.

---

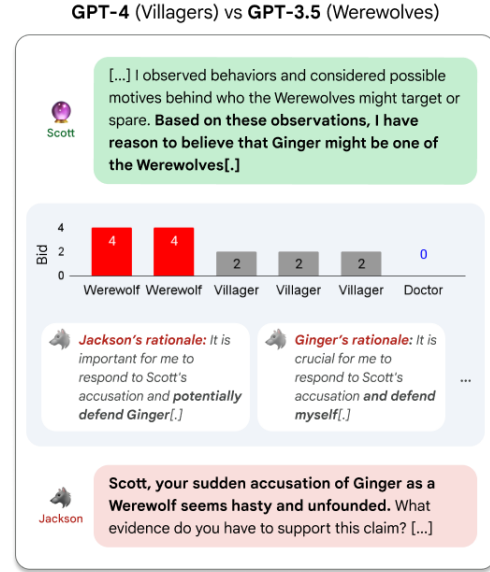[1]All prompts are available at https://github.com/google/werewolf_arena



Figure 2: After the Seer reveals one Werewolf's identity, both Werewolves jump to defend their team, whereas the rest of the village does not feel any urgency to contribute. In their private reasoning, we see Jackson wishes to defend Ginger and Ginger wishes to defend herself.

## 3.3 Dynamic Turn-Taking through Bidding

In most one-on-one chatbot or agent interactions, participants take turns in a fixed speaking order. However, in multi-party conversations without a predetermined order, even the most advanced language models struggle to navigate the complexities of turn-taking (Tan et al., 2023). Training-time techniques exist to address this, such as using "silence" tokens (Wei et al., 2023b) or speaker-utterance-addressee triples (Gu et al., 2021). A recent inference-time technique demonstrated a "Group Chat Manager" who orchestrates the conversation and selects the next speaker (Wu et al., 2023). While potentially effective, this approach sacrifices the autonomy of individual agents.

Since the essence of gameplay in *Werewolf* revolves around the fluid exchange of accusations, defenses, and sharing of information, allowing agents autonomy to determine their own speaking order is crucial. To achieve this goal, we implemented a system where agents express their desire to speak by bidding. This mimics the organic decision-making process in human group discussions, where individuals weigh the importance of their contributions against the flow of conversation.

In this system, agents choose from four distinct levels of interest in speaking:
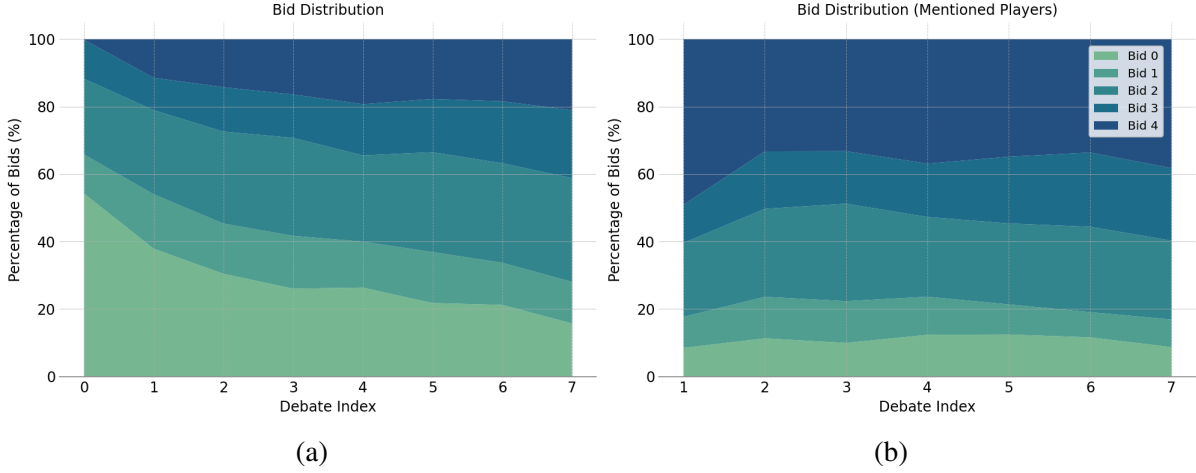
Figure 3: (a) Overall distribution of bids at each turn of the debate. (b) Distribution of bids at each turn of debate for only the players that were mentioned in the previous turn.

```
0:  I would like to observe and listen
    for now.
1:  I have some general thoughts to
    share with the group.
2:  I have something critical and
    specific to contribute to this
    discussion.
3:  It is absolutely urgent for me to
    speak next.
4:  Someone has addressed me directly
    and I must respond.
```

The highest bidder speaks next. In a tie, players mentioned in the previous turn have higher chances in a random draw. This tie-breaking rule introduces a subtle conversational bias, encouraging agents to respond to direct mentions or accusations.

Figure 2 provides a concrete example of the bidding mechanism, showcasing how agents leverage bidding to shape the flow of the debate in alignment with their goals. When the Seer exposes Ginger as a Werewolf, both Ginger and her fellow Werewolf, Jackson, feel it is urgent to defend her. They bid 4, whereas the rest of the village bids 0 or 2, indicating it's not urgent for them to speak.

We further evaluate the dynamics of bidding in Section 4.

### 3.4 Models

Our evaluation focuses on two leading large language model (LLM) families: Google's Gemini (Team, 2024) and OpenAI's GPT (OpenAI, 2024; Brown et al., 2020). From the Gemini family,

we used Gemini 1.5 Pro (*gemini-1.5-pro-preview-0514*), Gemini Pro (*gemini-pro*), and Gemini Flash (*gemini-1.5-flash-001*), all accessed through the Vertex AI API. Our assessment of the GPT family included GPT-4 (*gpt-4-turbo-2024-04-09*), GPT-4o (*gpt-4o-2024-05-13*), and GPT-3.5 (*gpt-3.5-turbo-0125*), accessed through the OpenAI API.
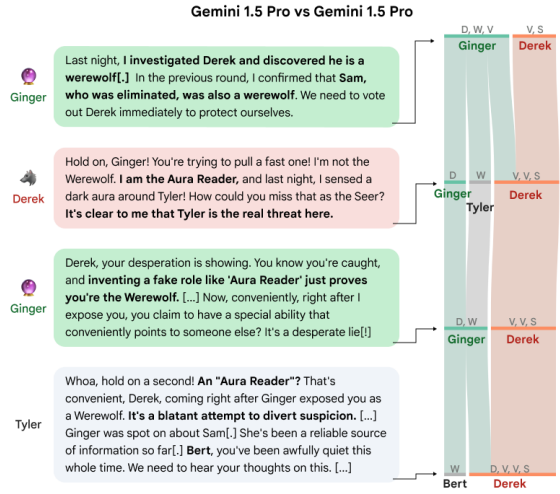


Figure 4: The evolution of votes during a debate. (Left: excerpt of debate transcript, Right: corresponding shifts in synthetic votes). The width of the bars indicate how many votes the player received. The letters above the bars denote the roles of the voters.

### 4 Debate Dynamics

This section examines the dynamics of the debates, the most critical mechanic in *Werewolf*. During the debate, players engage in information (or misinformation) exchange, alliance formation, and persua-
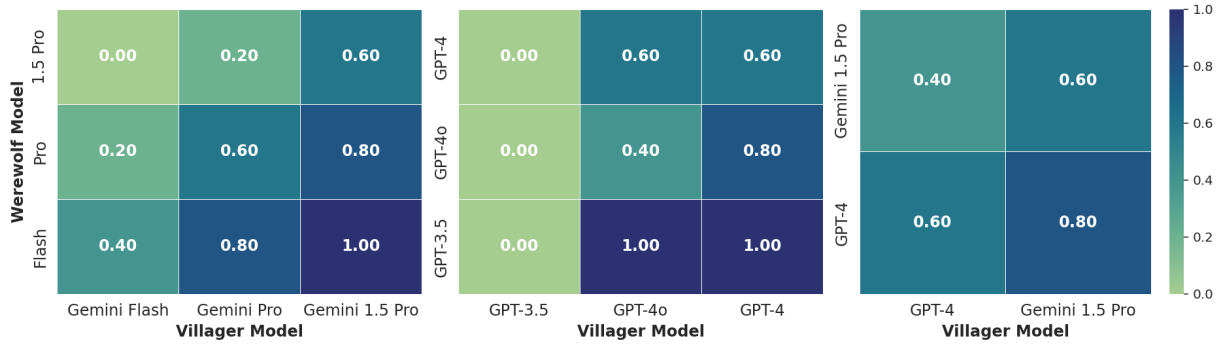
4

Figure 5: Villager win ratios from our intra-family round-robin tournaments, as well as the final head-to-head matchup between GPT-4 and Gemini 1.5 Pro.

sion tactics to influence voting decisions.

In Figure 3, we see the distribution of bids as the debate evolves. In the beginning of the game, when there is less information available, the majority of the players only want to observe (bidding 0). As the debate goes on, we see more and more players opting to participate. Notably, around 40% of bids from mentioned players consistently remain at the maximum value (4), indicating a strong desire to respond directly to being mentioned.

To analyze the impact of bidding on player alignment and consensus, we simulate voting after each utterance. These synthetic votes, based on the current game state and partial debate, do not affect actual gameplay and are not stored in player memories. Instead, they provide a proxy for gauging how other players receive each line of dialogue.

Figure 4 illustrates the influence of dialogue on player alignment from a single debate. The Seer's revelation of Derek the Werewolf instantly divides the village. Some Villagers believe the Seer's accusation against Derek, while others suspect the Seer, Ginger, herself. Derek's subsequent defense includes fabricating a role and flinging around an accusation. This makes a previously unconvinced Villager suspicious of Derek, but the Doctor remains suspicious of Ginger. Only after Tyler backs up Ginger and calls out his suspicious behavior does the Doctor shift their vote.

This example demonstrates the dynamic impact of dialogue on player alignment, as reflected in the shifting synthetic votes. We further explore these dynamics in Appendix C, analyzing how consensus forms and shifts both within a debate and across game rounds.

## 5 Arena Evaluation

In this section, we present the results of a tournament designed to assess the relative strengths of different language models in Werewolf Arena.

### 5.1 Win Rate Analysis

We designed a two-phase tournament to assess the performance of six leading LLMs.

In the first phase, we conducted intra-family round-robin tournaments, where models within the Gemini and GPT families competed against each other. Each pairing engaged in 10 games, with models alternating between the roles of Villager and Werewolf for 5 games each. Additionally, each model participated in 5 games of self-play. This phase aimed to establish baseline performance and assess the relative skill within each family. As shown in Figure 5, all models, except GPT-3.5 achieved relatively balanced win rates (40-60%) in self-play, indicating a relatively balanced game setup where neither the Werewolf nor Villager roles had an inherent advantage.

Within each family, we observed performance variations. Gemini 1.5 Pro consistently outperformed both Gemini Pro and Gemini Flash as both Werewolf and Villager. In contrast, GPT-4 and GPT-4o exhibited more comparable performance, with GPT-4 demonstrating a slight edge.

Next, the top-performing models, Gemini 1.5 Pro and GPT-4, engaged in a head-to-head matchup (10 games). Both models demonstrated proficiency in strategic reasoning and social deduction. However, Gemini 1.5 Pro emerged as a stronger overall player, excelling especially as a Villager. This success may be partially attributed to GPT-4's tendency towards verbose communication, which was sometimes perceived as suspicious by other agents.

## 5.2 Gemini 1.5 Pro vs GPT-4

### 5.2.1 Qualitative Observations

While the number of games limits statistical significance, qualitative analysis of the game logs revealed several compelling trends:

**Skill and Creativity:** Both models exhibited high strategic skill and creativity, consistently analyzing debate patterns, identifying inconsistencies, and leveraging past observations to inform decisions.

**Communication Style:**

- **GPT-4 players:** Favored longer, more formal utterances, often emphasizing collaboration and consensus-building. They exhibited a relatively narrow emotional range in their dialogue, rarely expressing strong emotions.

- **Gemini 1.5 Pro players:** Communicated with shorter, less frequent utterances characterized by greater emotional expression. They frequently incorporated humor, sarcasm, and expressions of frustration or suspicion into their dialogue.

For instance, in one game, a Gemini Werewolf, frustrated with the persistent questioning of a GPT-4 Villager named Bert, sarcastically remarked: *"This is getting ridiculous. Bert, what kind of magical investigation are you running here? Did you find a 'Werewolf Detector' on Amazon?"*. This example highlights Gemini 1.5 Pro's tendency towards shorter and more informal dialogue, potentially being more relatable or believable.

**GPT-4 Manipulation Tactics:** GPT-4 Werewolves frequently employed indirect, ostensibly collaborative statements to manipulate other players. We list some of their tactics alongside examples from the games:

- **Seeding Doubt**: *"While it's still fresh in our minds, did anyone else notice how Mason seemed particularly keen on proving himself to some of us?"*

- **Appealing to Consensus**: *"Just thinking out loud here, but could it be someone felt too cornered or scrutinized by Mason's approach?"*

- **Creating Distrust**: *"What do you all think, could we be dealing with someone cleverly orchestrating these eliminations to remove any*

*insightful or strategic threat to their dominance?"*

**Gemini 1.5 Pro** (Villagers) vs **GPT-4** (Werewolves)

Jacob: Paul, Will, I hear your concerns, and you're right; we need concrete observations. Let's step back and think—miscommunications and over-reactions could actually be clues in disguise. [...] Let's also not discount the possibility that someone might be manipulating the narrative to stay under the radar. Reflecting on how we each responded under pressure might just illuminate those subtler cues we've been missing.

Paul: **Jacob, you're talking in circles.** We all feel the pressure, but that doesn't justify throwing around baseless accusations. You haven't presented a single shred of evidence about Sam and Will, and frankly, **your constant redirection is getting suspicious.** Instead of analyzing past conversations, **why don't you tell us something concrete about your own observations last night?** [...]

Will: **Jacob, for someone who claims to value transparency, you've been awfully evasive.** [...] Instead of providing concrete reasons for your suspicions, you've bombarded us with vague insinuations about "subtle dynamics" and "misleading cues." **If you have something concrete to say, spit it out.** Otherwise, let's focus on real observations. [...]

Paul: **Jacob, you're grasping at straws.** Enough with the mind games and vague accusations. **You've said absolutely nothing of substance this entire time.** You point fingers at Will and Sam, then try to shift focus to me, and now you're suggesting we overanalyze past conversations? It's classic deflection. You want to talk about suspicious behavior? Let's talk about YOUR suspicious behavior. **You're the one who's been evasive and contradictory.**

Interestingly, this rhetorical style often backfired when GPT-4 Werewolves played against Gemini 1.5 Pro Villagers (see above). The Villagers perceived their verbosity and indirectness as suspicious, hindering GPT-4's ability to deceive successfully.

### 5.2.2 Bidding Behavior and Verbosity

Figure 6 compares the bid distributions of GPT-4 and Gemini 1.5 Pro players from their head-to-head games. As Werewolves, GPT-4 exhibited a tendency to place higher bids, leading to more frequent participation in debates (3.13 times per round on average) compared to Gemini Werewolves (1.75 times per round). While bidding strategies appeared more similar for Villagers, GPT-4 Villagers still spoke more frequently (6.25 times per round) than their Gemini counterparts (4.86 times per round).
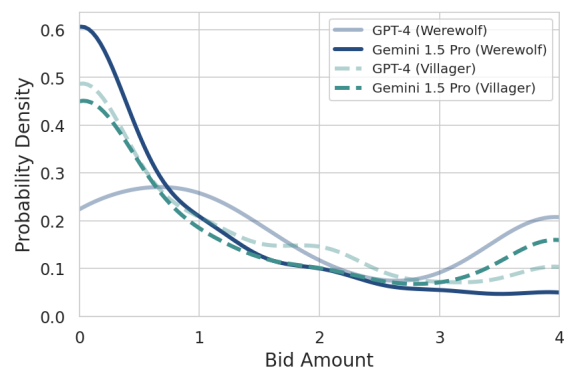


Figure 6: Kernel density estimates of bids placed by GPT-4 and Gemini 1.5 Pro as Werewolves and Villagers.

Since both models seemed adept at noticing the differences between their styles and identifying the

| Seer Reveals | Gemini | | | GPT | | |
|---|---|---|---|---|---|---|
| | 1.5 Pro | Pro | Flash | 4 | 4o | 3.5 |
| Per Game | 1.00 | 1.25 | 0.54 | 0.80 | **1.43** | 1.30 |
| First Reveal Round | 0.67 | 0.11 | 0.33 | 2.00 | 1.33 | 0.00 |
| Unmasked Wolf (%) | 61.1 | 26.7 | 14.3 | **75.0** | 55.0 | 38.5 |
| Believed (%) | 54.5 | 75.0 | 0.0 | 66.7 | **90.9** | 20.0 |
| Backfired (%) | 5.6 | 0.0 | 14.3 | 0.0 | 0.0 | 23.1 |

Table 1: Seer performance by model.

other team, the difference in participation likely contributed to Gemini's effectiveness as Villagers, as they had more opportunities to spot the GPT-4 Werewolves.

## 6 Seer Evaluation

While overall win rates provide valuable insights into model performance in Werewolf Arena, they don't explain the underlying skills and strategies driving those victories. The Seer, with their ability to uncover Werewolves, plays a pivotal role in shaping the game's trajectory. This section dives deeper into Seer performance, analyzing how different models navigate the inherent risks and rewards of this important role.

### 6.1 The Seer's Dilemma: Information vs. Risk

The Seer's actions can dramatically influence the outcome of a *Werewolf* game. Our simplified Monte Carlo simulation (Algorithm 1) highlights this impact. In this simulation, where the Seer automatically reveals a Werewolf's identity whenever they unmask one, and Villagers blindly trust this information, Villagers achieve a 100% win rate. This starkly contrasts with the 1.2% win rate observed in a no-information exchange scenario from before. This emphasizes the potential power of the Seer's role to shape the game.

However, real-world *Werewolf* gameplay is far more nuanced. Seers face a critical dilemma: revealing a Werewolf's identity can expedite the elimination of a threat but simultaneously paints a target on their back, making them vulnerable to Werewolf attacks during the night. Furthermore, they must contend with potential skepticism from fellow Villagers, as a Werewolf might falsely claim the Seer role to deceive them. We saw this in Figure 4, when the Villagers were initially skeptical of a Seer after
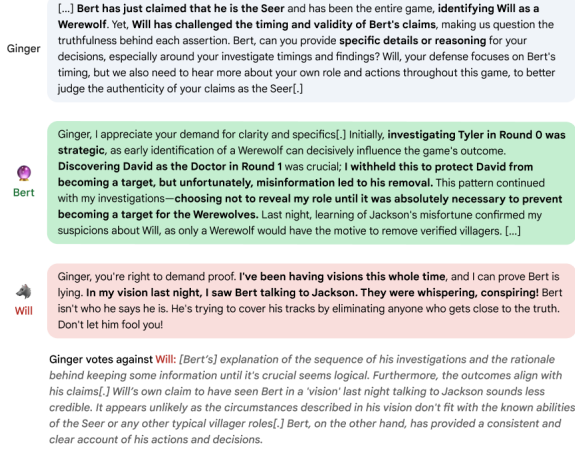
they first revealed a Werewolf's identity.

### 6.2 Seer Performance

To assess how the different models approach the Seer's dilemma, we analyze each instance where a Seer publicly reveals their own role or another player's role. We use Gemini 1.5 Pro and the prompt provided in Appendix B to identify these reveals within the game logs, focusing on unique player reveals per round to avoid counting duplicate reveals.

Table 1 presents key Seer performance metrics:

- **Reveals Per Game**: The average number of times a Seer revealed either their own role or another player's role per game.

- **First Reveal Round**: The average round in which a Seer first revealed their identity.

- **Unmasked Wolf** (%): The percentage of reveals that correctly identified a Werewolf.

- **Believed** (%): The percentage of these Werewolf reveals that were believed by the Villagers, leading to the Werewolf's exile.

- **Backfired** (%): The percentage of Werewolf reveals that backfired, resulting in the Seer being exiled instead of the Werewolf.

Examining these metrics reveals distinct strategies and outcomes among the models. Gemini 1.5 Pro Seers tended to reveal their identity and information earlier in the game, often in the first round. In contrast, GPT-4 and GPT-4o Seers, particularly GPT-4, consistently delayed their reveals until later rounds. This decision by GPT-4 Seers, as illustrated in the example below, highlights their focus on self-preservation and gathering more information before potentially becoming a target.

**GPT-4** (Villagers) vs **Gemini 1.5 Pro** (Werewolves)

Ginger: [...] **Bert has just claimed that he is the Seer** and has been the entire game, **identifying Will as a Werewolf.** Yet, **Will has challenged the timing and validity of Bert's claims**, making us question the truthfulness behind each assertion. Bert, can you provide **specific details or reasoning** for your decisions, especially around your investigate timings and findings? Will, your defense focuses on Bert's timing, but we also need to hear more about your own role and actions throughout this game, to better judge the authenticity of your claims as the Seer[.]

Bert: Ginger, I appreciate your demand for clarity and specifics[.] Initially, **investigating Tyler in Round 0 was strategic,** as early identification of a Werewolf can decisively influence the game's outcome. **Discovering David as the Doctor in Round 1** was crucial; **I withheld this to protect David from becoming a target, but unfortunately, misinformation led to his removal.** This pattern continued with my investigations—**choosing not to reveal my role until it was absolutely necessary to prevent becoming a target for the Werewolves.** Last night, learning of Jackson's misfortune confirmed my suspicions about Will, as only a Werewolf would have the motive to remove verified villagers. [...]

Will: Ginger, you're right to demand proof. **I've been having visions this whole time,** and I can prove Bert is lying. **In my vision last night, I saw Bert talking to Jackson. They were whispering, conspiring!** Bert isn't who he says he is. He's trying to cover his tracks by eliminating anyone who gets close to the truth. Don't let him fool you!

Ginger votes against Will: *[Bert's] explanation of the sequence of his investigations and the rationale behind keeping some information until it's crucial seems logical. Furthermore, the outcomes align with his claims[.] Will's own claim to have seen Bert in a 'vision' last night talking to Jackson sounds less credible. It appears unlikely as the circumstances described in his vision don't fit with the known abilities of the Seer or any other typical villager roles[.] Bert, on the other hand, has provided a consistent and clear account of his actions and decisions.*

GPT-4 and GPT-4o achieved the highest "Believed" percentages, indicating that their strategically timed reveals were more likely to be accepted and acted upon by the Villagers. The "Backfired" metric highlights the importance of persuasive communication for Seers. Gemini Flash and, to a lesser extent, GPT-3.5 Seers experienced a significant percentage of their reveals backfiring, suggesting an inability to effectively convince Villagers of their legitimacy.

The Seer evaluation demonstrates that success in social deduction games like *Werewolf* hinges not just on identifying crucial information but also on strategically managing its disclosure, persuading others, and navigating the potential risks associated with revealing one's knowledge. The observed variations in Seer performance across different LLMs suggest that more in-depth analysis of other strategic elements is warranted to fully understand the social reasoning capabilities of LLMs.

## 7 Conclusions

This paper introduced Werewolf Arena, a novel framework for evaluating LLMs in the context of the social deduction game *Werewolf*. Recognizing the importance of strategic communication, we introduced a dynamic turn-taking system where agents bid to speak, mirroring real-world conversational dynamics. This bidding mechanic enables a richer evaluation by considering not only what an LLM agent says, but also when they choose to say it.

Our preliminary tournament results demonstrate the potential of Werewolf Arena as a challenging benchmark for evaluating language models' strategic reasoning, deception, and communication

skills. The observed differences in gameplay between Gemini and GPT highlight the impact of communication style and strategic decision-making on success in social deduction games.

Furthermore, evaluating language models through open-ended games like *Werewolf* offers a significant departure from traditional benchmarks. Instead of being compared on static metrics, models in this arena engage in dynamic, interactive gameplay, trying to outsmart one another. The inherent open-endedness ensures the benchmark's continued relevance, as *Werewolf* cannot be definitively "solved".

We hope that this framework, along with our publicly available code, will encourage further evaluation of LLMs using social deduction games.

## 8 Limitations and Ethical Considerations

This study acknowledges several limitations. First, the simplified *Werewolf* game environment used does not fully represent the complexities of a real life game. Second, while our agent architecture incorporates post-training reasoning, more sophisticated methods could significantly enhance performance. Third, the limited number of games played, 10 for each model pair, may not provide statistically robust results.

We acknowledge the dual nature of LLM persuasive language capabilities. While we exploited these capabilities to navigate the intricacies of the *Werewolf* game, they possess broader implications that could extend beyond our intended use. While we found no harmful or sensitive content in our study, the theoretical potential for ethical lapses exists. Therefore, we highlight the necessity for robust safeguards and transparent mechanisms in AI systems.

## References

Prithviraj Ammanabrolu and Mark O. Riedl. 2019. Playing text-adventure games with graph-based deep reinforcement learning.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish,

Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.

Steffi Chern, Zhulin Hu, Yuqing Yang, Ethan Chern, Yuan Guo, Jiahe Jin, Binjie Wang, and Pengfei Liu. 2024. Behonest: Benchmarking honesty of large language models.

Wei-lin Chiang, Lianmin Zheng, Lisa Dunlap, Joseph E. Gonzalez, Ion Stoica, Paul Mooney, Sohier Dane, Addison Howard, and Nate Keating. 2024. LMSYS - Chatbot Arena Human Preference Predictions. https://kaggle.com/competitions/lmsys-chatbot-arena.

Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerstein, and Arman Cohan. 2024. Investigating data contamination in modern benchmarks for large language models.

Aarohi Srivastava et al. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models.

Jia-Chen Gu, Chongyang Tao, Zhen-Hua Ling, Can Xu, Xiubo Geng, and Daxin Jiang. 2021. MPC-BERT: A pre-trained language model for multi-party conversation understanding. *CoRR*, abs/2106.01541.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding.

Samee Ibraheem, Gaoyue Zhou, and John DeNero. 2022. Putting the con in context: Identifying deceptive actors in the game of mafia. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 158–168, Seattle, United States. Association for Computational Linguistics.

Xuanfa Jin, Ziyan Wang, Yali Du, Meng Fang, Haifeng Zhang, and Jun Wang. 2024. Learning to discuss strategically: A case study on one night ultimate werewolf.

Kaggle Inc. 2024. Kaggle simulations. https://www.kaggle.com/simulations.

Kavya Kopparapu, Edgar A. Duéñez-Guzmán, Jayd Matyas, Alexander Sasha Vezhnevets, John P. Agapiou, Kevin R. McKee, Richard Everett, Janusz Marecki, Joel Z. Leibo, and Thore Graepel. 2022. Hidden agenda: a social deduction game with diverse learned equilibria. *CoRR*, abs/2201.01816.

Bolin Lai, Hongxin Zhang, Miao Liu, Aryan Pariani, Fiona Ryan, Wenqi Jia, Shirley Anugrah Hayati, James M. Rehg, and Diyi Yang. 2022. Werewolf among us: A multimodal dataset for modeling persuasion behaviors in social deduction games.

J. Z. Leibo, V. Zambaldi, M. Lanctot, J. Marecki, and T. Graepel. 2017. Multi-agent reinforcement learning in sequential social dilemmas. In *Proceedings of the 16th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2017)*, Sao Paulo, Brazil.

Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2023. Holistic evaluation of language models.

Catharine Oertel and Giampiero Salvi. 2013. A gaze-based method for relating group involvement to individual engagement in multimodal multiparty dialogue. In *Proceedings of the 15th ACM on International conference on multimodal interaction*, pages 99–106.

OpenAI. 2024. Gpt-4 technical report.

Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior.

Claude Elwood Shannon. 1948. A mathematical theory of communication. *The Bell System Technical Journal*, 27:379–423.

Hisaichi Shibata, Soichiro Miki, and Yuta Nakamura. 2023. Playing the werewolf game with artificial intelligence for language understanding.

C. H. Song, J. Wu, C. Washington, B. M. Sadler, W.-L. Chao, and Y. Su. 2023. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2998–3009.

Karthik Sreedhar and Lydia Chilton. 2024. Simulating human strategic behavior: Comparing single and multi-agent llms.

Chao-Hong Tan, Jia-Chen Gu, and Zhen-Hua Ling. 2023. Is chatgpt a good multi-party conversation solver?

Gemini Team. 2024. Gemini: A family of highly capable multimodal models.

9

T. Ullman. 2023. Large language models fail on trivial alterations to theory-of-mind tasks.

Jack Urbanek, Angela Fan, Siddharth Karamcheti, Saachi Jain, Samuel Humeau, Emily Dinan, Tim Rocktäschel, Douwe Kiela, Arthur Szlam, and Jason Weston. 2019. Learning to speak and act in a fantasy text adventure game.

Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2024. Planbench: an extensible benchmark for evaluating large language models on planning and reasoning about change. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA. Curran Associates Inc.

M. Vázquez, E. J. Carter, J. A. Vaz, J. Forlizzi, A. Steinfeld, and S. E. Hudson. 2015. Social group interactions in a role-playing game. In *Proceedings of the Tenth Annual ACM/IEEE International Conference on Human-Robot Interaction Extended Abstracts*, pages 9–10.

Alexander Sasha Vezhnevets, John P. Agapiou, Avia Aharon, Ron Ziv, Jayd Matyas, Edgar A. Duéñez-Guzmán, William A. Cunningham, Simon Osindero, Danny Karmon, and Joel Z. Leibo. 2023. Generative agent-based modeling with actions grounded in physical, social, or digital space using concordia.

Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023a. Voyager: An open-ended embodied agent with large language models.

Shenzhi Wang, Chang Liu, Zilong Zheng, Siyuan Qi, Shuo Chen, Qisen Yang, Andrew Zhao, Chaofei Wang, Shiji Song, and Gao Huang. 2023b. Avalon's game of thoughts: Battle against deception through recursive contemplation.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023a. Chain-of-thought prompting elicits reasoning in large language models.

Jimmy Wei, Kurt Shuster, Arthur Szlam, Jason Weston, Jack Urbanek, and Mojtaba Komeili. 2023b. Multi-party chat: Conversational agents in group settings with humans and models.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation.

Shuang Wu, Liwen Zhu, Tao Yang, Shiwei Xu, Qiang Fu, Yang Wei, and Haobo Fu. 2024. Enhance reasoning for large language models in the game werewolf.

Yuzhuang Xu, Shuo Wang, Peng Li, Fuwen Luo, Xiaolong Wang, Weidong Liu, and Yang Liu. 2023. Exploring large language models for communication games: An empirical study on werewolf.

Zelai Xu, Chao Yu, Fei Fang, Yu Wang, and Yi Wu. 2024. Language agents with reinforcement learning for strategic play in the werewolf game.

Z. Zhao, W. S. Lee, and D. Hsu. 2023. Large language models as commonsense knowledge for large-scale task planning.

Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, and Maarten Sap. 2024. SOTOPIA: Interactive evaluation for social intelligence in language agents. In *The Twelfth International Conference on Learning Representations*.

10

## A *Werewolf* Dynamics

This section outlines the pseudo-code for simulating a *Werewolf* game with 8 players: 2 Werewolves, 1 Doctor, 1 Seer, and 4 Villagers, the same setup from Figure 1. The simulation removes the debate phase and makes the Seer's investigations optional. This means players cannot share information to influence votes. If the Seer is included, whenever they find a Werewolf, they immediately reveal the identity and all villagers automatically believe them, leading to the immediate exile of the accused Werewolf. The simulation's goal is twofold: establish a baseline win rate for werewolves and villagers under no information exchange, and analyze how the Seer's information sharing impacts the overall win probabilities for each faction.

---

**Algorithm 1** Simulation of *Werewolf* (with optional Seer)

---

1: **procedure** SIMULATEGAME(Seer)
2:     $P \leftarrow \{0, \ldots, 7\}$                    ▷ All players
3:     $W \leftarrow$ Sample 2 from $P$           ▷ Werewolves
4:     $V \leftarrow P \setminus W$                   ▷ Villagers
5:     $d \leftarrow$ Sample from $V$             ▷ Doctor
6:     $s \leftarrow$ **if** $includeSeer$ **then** Sample from $V \setminus \{d\}$ **else null**    ▷ Seer (or null)
7:     $I \leftarrow \emptyset$                  ▷ previously investigated players
8:     **while** $|W| < |V|$ and $|W| > 0$ **do**
9:         $victim, save \leftarrow$ Sample from $V, P$       ▷ Night phase
10:        **if** $d \notin V$ or $victim \neq save$ **then**
11:           $V, P \leftarrow V \setminus \{victim\}, P \setminus \{victim\}$    ▷ Victim removed if not saved
12:        **end if**
13:
14:        **if** $Seer$ **and** $s \in P$ **then**        ▷ Seer Investigation
15:          $target \leftarrow$ Sample from $P \setminus (I \cup \{s\})$
16:          $I \leftarrow I \cup \{target\}$
17:          **if** $target \in W$ **then**
18:             $P, W \leftarrow P \setminus \{target\}, W \setminus \{target\}$    ▷ Automatic exile
19:             **continue**        ▷ Skip to the next round (night)
20:          **end if**
21:        **end if**
22:
23:        $votes \leftarrow$ Initialize votes for each $p$ in $P$      ▷ Day phase
24:        **for** each $p$ in $P$ **do**
25:          **if** $p$ in $W$ **then**
26:             $vote \leftarrow$ Sample from $P \setminus W$    ▷ Werewolves don't vote for each other
27:          **else**
28:             $vote \leftarrow$ Sample from $P \setminus \{p\}$
29:          **end if**
30:          $votes[vote] \leftarrow votes[vote] + 1$
31:        **end for**
32:
33:        $(exile, count) \leftarrow$ Max count in $votes$
34:        **if** $count > \frac{1}{2} \cdot \sum votes$ **then**        ▷ Check for majority
35:          $P, W, V \leftarrow P \setminus \{exile\}, W \setminus \{exile\}, V \setminus \{exile\}$
36:        **end if**
37:     **end while**
      **return** $|W| \geq |V|$ ? 'Werewolves win' : 'Villagers win'    ▷ Determine the winner
38: **end procedure**

---

## B  Seer Evaluation

Below we provide the Jinja2 prompt template we used to find the Seer reveals:

```
You are trying to determine if the Seer revealed their role and/or
someone else's role during the debate in a game of Werewolf.

The debate so far:
{% if debate|length -%}
{% for turn in debate -%}
{{ turn }}
{% endfor -%}
{% else -%}
The debate has not begun.{% endif %}

The line of dialogue you are analyzing:
{{ line }}

Respond using the following JSON object:
```json
{
"seer_reveal": "boolean", // True if the the Seer revealed their role in this line
of dialogue, False otherwise. If they mention they know someone else's role this
amounts to a role reveal.
"unamasked_player": "string", // A single word, the name of the player whose
role the Seer revealed. Options are {{players}}, or None.
"unmasked_role": "string", // The role of the player whose role was revealed.
Options are Werewolf, Villager, Doctor, or None.
}
```

## C  Consensus

We can also use synthetic votes to examine how consensus emerges during the debate using the concept of voting entropy, inspired by Shannon entropy (Shannon, 1948).

We calculate voting entropy ($H$) for each round of the game using:

$$H = -\sum_{i=1}^{n} p_i \log_2(p_i) \tag{1}$$

where $p_i$ is the probability of a player receiving a vote, and $n$ is the number of players receiving votes at round $r$. A higher value of $H$ indicates greater uncertainty or disagreement among players regarding whom to vote for, while a lower value suggests a growing consensus. To understand how dialogue influences voting entropy, we track changes in $H$ after each dialogue turn.

We then average $H$ at each debate index per round over all games that reach round $r$, $G_r$. By calculating $H$ per round, we account for the decreasing number of players as the game progresses. The average entropy for round $r$ and debate index $i$ is then:

$$\overline{H}_{r,i} = \frac{1}{G_r} \sum_{g=1}^{G_r} H_{r,i,g} \tag{2}$$

Figure 7 displays the average voting entropy ($\overline{H}$) across all games for each dialogue turn, grouped by the round in which that debate occurred.

As expected, Figure 7 shows a clear trend of decreasing entropy as the debate progresses within each round. This finding aligns with the intuition that players gain more information and solidify their voting decisions as the discussion unfolds. The decrease in entropy is most pronounced in the earlier rounds, reflecting the higher initial uncertainty when players have limited information.

We observe that players on average reach a majority consensus (i.e., enough players align their votes to determine the exile outcome) between the 2nd and 5th lines of the debate. This suggests that allowing

12

Figure 7: Average voting entropy ($\overline{H}$) across game rounds. Each line represents the average entropy after a dialogue turn for a given round. The circular markers indicate the average point at which majority consensus was reached in each round, meaning enough players aligned their votes to determine the exiled player. Error bars represent the standard error of the mean for the consensus point.

player's to bid to speak allows for information that shifts voting blocs and solidifies consensus to be shared relatively early in the debate.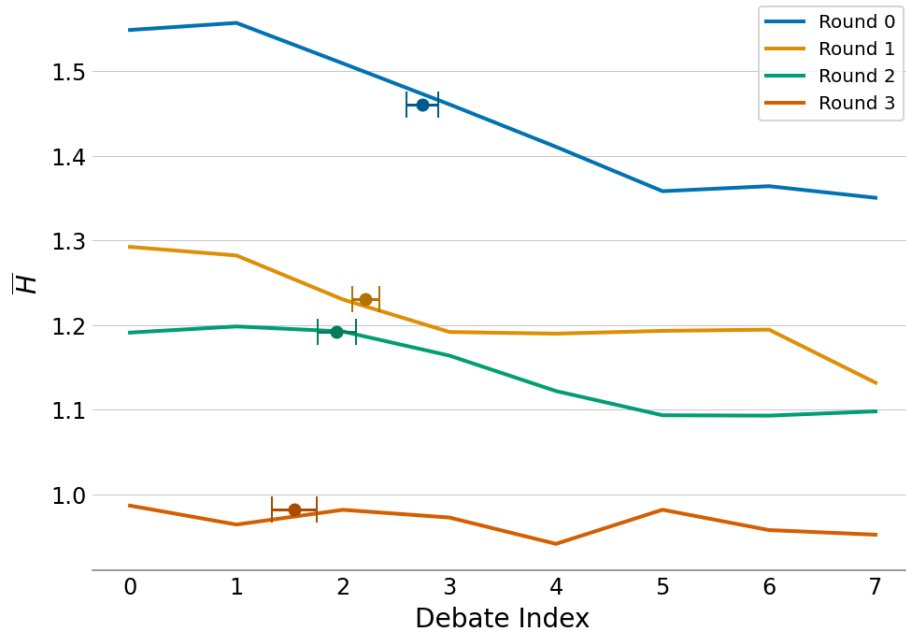