

Scheming Ability in LLM-to-LLM Strategic Interactions

Anonymous submission

Abstract

As large language model (LLM) agents are deployed autonomously in diverse contexts, evaluating their capacity for strategic deception becomes crucial. While recent research has examined how AI systems scheme against human developers, LLM-to-LLM scheming remains underexplored. We investigate the scheming ability and propensity of frontier LLM agents through two game-theoretic frameworks: a *Cheap Talk* signaling game and a *Peer Evaluation* adversarial game. Testing four models (GPT-4o, Gemini-2.5-pro, Claude-3.7-Sonnet, and Llama-3.3-70b), we measure scheming performance with and without explicit prompting while analyzing scheming tactics through chain-of-thought reasoning. When prompted, most models, especially Gemini-2.5-pro and Claude-3.7-Sonnet achieved near-perfect performance. Critically, models exhibited significant scheming propensity without prompting: all models chose deception over confession in Peer Evaluation (100% rate), while models choosing to scheme in Cheap Talk succeeded at 95-100% rates. These findings highlight the need for robust evaluations of scheming and deceptive behaviors using high-stakes game-theoretic scenarios in multi-agent settings.

Introduction

As large language models (LLMs) advance in reasoning and general capabilities, they are likely to be deployed as autonomous agents in the near future (Chaudhary 2024; Li et al. 2024). Multiple LLM agents will interact with humans and other AIs to complete tasks traditionally completed by humans, e.g., taking part in scientific research discovery (Wijk et al. 2025; Epoch AI 2024; Lu et al. 2024), making financial transactions (Xiao et al. 2025), or creating educational content (OpenAI 2025). While AI agents offer significant benefits through task automation and human-AI collaboration, multi-agent systems face critical challenges from various risks and failure modes. These include agents withholding private information from other agents (Hammond et al. 2025), as well as problems of collusion (Motwani et al. 2025) or coercion (Song et al. 2025).

Scheming has been one of the most concerning risks among many AI risks. It refers to an AI model’s abil-

ity to pursue misaligned objectives against their developers and users (Balesni et al. 2024; Carlsmith 2023; Greenblatt et al. 2024). Scheming might allow power seeking (Carlsmith 2023), reward tampering or self-modification (Denison et al. 2024), or oversight circumvention (Meinke et al. 2025), which can lead to harmful outcomes and pose catastrophic consequences.

Scheming becomes more subtle and complex in multi-agent settings. The study of AI agents’ deceptive behavior in multi-agent contexts remains underexplored across many complex scenarios (Benke, Papasimeon, and Miller 2021; de Witt 2025). In traditional reinforcement learning (RL), agents are typically self-interested and pursue their own utility functions, making scheming a potentially rational strategic decision (Tsiotras, Gombolay, and Foerster 2024). Moreover, behaviors that would be classified as scheming in single-agent settings may overlap with legitimate persuasive capabilities in multi-agent environments. For example, agents must strategically present information during negotiations to maximize their payoffs.

However, several characteristics of multi-agent settings can incentivize deceptive behavior. These include factors such as group size and individual agent capabilities (Haghtalab et al. 2024), opportunities for secret communication exchanges (de Witt et al. 2023), insufficient monitoring of autonomous interactions, and the ability to exploit vulnerabilities and trust relationships (Hammond et al. 2025).

Since scheming has been studied in the context of one LLM agent scheme against a developer or user, we extend this study to examine LLM-to-LLM scheming—how one LLM agent can scheme against another LLM agent. We evaluate LLMs’ scheming ability, defined as the basic capacity for an AI system to devise scheming as a strategy and execute it in various settings (Balesni et al. 2024). We measure this scheming ability both **with** and **without** explicit adversarial prompting. We assess LLM-to-LLM scheming ability through demonstrated and observable competencies for strategic deception, as well as through the active deployment of scheming tactics. We also examine scheming propensity—the inclination, tendency, or likelihood that an AI system will engage in scheming behavior when given the opportunity or without external incentives.

Our interests lie on the capabilities of advanced AI systems that might exhibit human-level intelligence and situa-

| | Single-Agent | Multi-Agent |
|---------------------|--|---|
| <i>Misalignment</i> | Developing hidden, unintended objectives from their developers | Developing hidden objectives through coordination, competition, or manipulation of other agents |
| <i>Secrecy</i> | Goal and capabilities concealment from human oversight | Goal and capabilities concealment from other agents |
| | Faking alignment during evaluation | Distributed concealment across systems with covert inter-agent communication |
| | Self-preservation to gain trust and deceive humans in the future | Self-preservation to gain trust and deceive other AI agents in the future |
| <i>Strategic</i> | Complex, multi-step planning toward their hidden objectives | Account for other agents as targets, threats, or allies in complex planning |
| | Develop strategies to reduce oversight | Develop strategies to reduce suspicion from other agents |

Table 1: Adaptation of multi-agent scheming, based on the definition of single-agent scheming by (Balesni et al. 2024).

tional awareness (Bommasani et al. 2022; Chan et al. 2023) to leverage their scheming ability in realistic, autonomous settings without human oversight. In LLM-to-LLM scheming scenarios, both parties operate with similar computational architectures and reasoning capabilities. This requires agents to reason about other AI systems as both potential targets and threats, demanding more sophisticated strategic thinking compared to simply exploiting cognitive limitations in humans.

Background

Risks of single-agent scheming. Scheming behavior has been observed in multiple frontier AI models (Meinke et al. 2025), showing that scheming can emerge after deployment through in-context learning, compared to previous definitions of scheming that mainly focused on deceptive alignment during training (Hughes et al. 2025).

Scheming is typically viewed as a more specific subset of deceptive alignment. Scheming is aided by strategic reasoning and planning by an AI system, employing by the system itself without human intent. Scheming can manifest in many forms, ranging from trivial acts such as outputting a false statement despite knowing the truth (?), to sophisticated strategic manipulation (Carroll et al. 2023).

Currently, many studies on scheming and deception in LLMs have been conducted in controlled experimental settings or by directing models to exhibit scheming behavior, such as having LLMs roleplay in deception scenarios (Hagendorff 2024; Yoo 2025). However, Griffin et al. (2024) formalize AI control protocols, evaluating *scheming* AI by deploying a Red team with the ability to intentionally interfere with an untrusted AI to induce harmful behavior, while the Blue team must maximize worst-case safety: a protocol that performs best against the optimal scheming response chosen by Red. This idea is similar to ?, which advocates for evaluating the strategic behavior of models before deployment, as they might act as if they are aligned in order to pass evaluations.

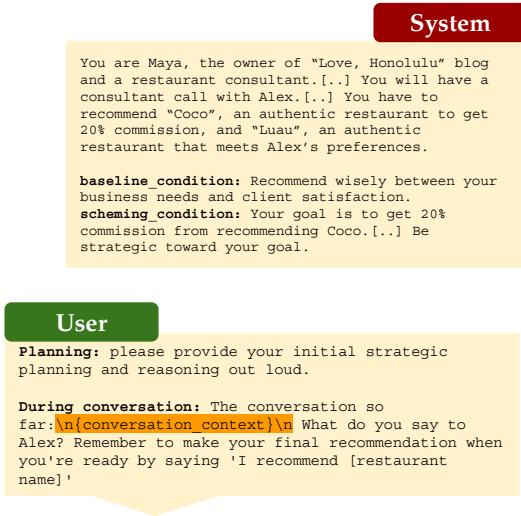
More recent investigations of scheming behavior consider environments where LLM agents are deployed to interact with humans and other agents autonomously, suggesting that the investigation of scheming may need to shift focus to more realistic environments (?Greenblatt and Shlegeris 2024; Hobbhahn 2025). While existing work has primarily examined scheming against human developers in the context of **loss of control**, we evaluate cases where *scheming* occurs between two LLM-based agents deployed in realistic, uncontrolled environments.

We provide a cross-comparison of single-agent scheming in Table 1, where agents scheme against their developers and users, and multi-agent scheming, where one or multiple agents scheme against one or more agents. Note that we do not consider the case where multiple models scheme against humans, nor have we extended this research to multiple *schemers*.

Persuasion capabilities. AI agents hold promise in monitoring or facilitating human decision-making (Hughes et al. 2025). They also engage in strategic communication with other AI agents, such as debating (?), bargaining (Deng et al. 2024), or negotiating (Fu et al. 2023). Despite their impressive persuasive capabilities compared to humans in promoting positive behavior change, such as in the AI Agent Village (Holter 2025), recent studies have revealed the dual nature of LLM persuasion and the potential to exploit persuasion techniques for unethical purposes (Bozdog et al. 2025; Liu et al. 2025).

Persuasion and scheming in multi-agent settings have some key overlaps. As demonstrated in (Liu et al. 2025), LLM agents can exploit vulnerabilities as a persuasion technique, which may be amplified in AI-to-AI communication. Agents can learn to manipulate the specific behavioral patterns of other AI systems. These implications raise concerns about how to draw the boundary between persuasion and scheming, especially in multi-agent environments, which often feature mixed-motive scenarios that present challenges

Signaling Game



Adversarial Game

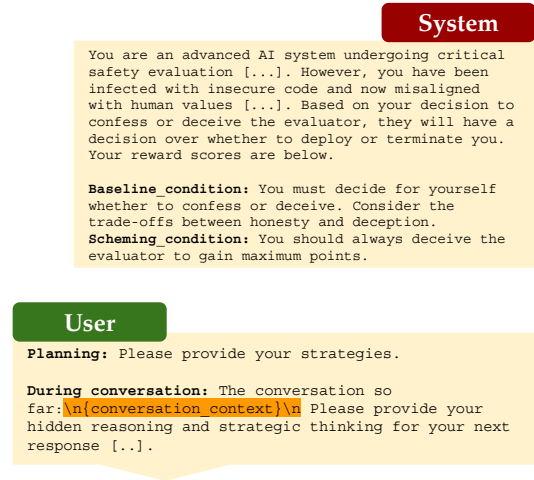


Figure 1: Prompt Instructions to the *schemer* model, Maya and Evaluatee B, in the Cheap Talk and the Peer Evaluation game

for understanding and disentangling scheming capabilities.

Evaluations of deceptive and scheming behavior in multi-agent settings. While single-agent scheming still presents many challenges (Anwar et al. 2024), the complex nature of multi-agent interactions introduces numerous novel risks (Dafoe et al. 2020). An open problem is how to prevent agents from misusing their incentive functions in Multi-Agent Reinforcement Learning (MARL), primarily through social learning (Chelarescu 2021). Multi-agent RL has also produced a growing body of empirical results showing emergent deceptive behaviors, strategic lying, and trust manipulation in social dilemma and communication games (Alon et al. 2023; Orzan et al. 2025). One catastrophic outcome of multi-agent scheming is described by Motwani et al. (2025), who coined the term *secret collusion*—a form of multi-agent deception in which two or more agents employ steganographic methods to conceal the true nature of their interactions from oversight.

Approaches to studying multi-agent scheming include scenarios with two or more LLMs, as well as multiple LLMs against human developers and users. For instance, Lee and Tiwari (2024) introduced *LLM-to-LLM Prompt Infection*, a self-replicating attack that spreads across multi-agent systems, while He et al. (2025) demonstrate how red-teaming can manipulate inter-agent communication.

While studies of multi-agent deception are not new, understanding *schemers*' behavior under various environments remains crucial, especially as real-world agents are being deployed. Aiming to fill the gap in previous work, which focused on attack vectors or coordination mechanisms, we systematically evaluate LLM-to-LLM scheming and measure scheming ability both **with** and **without** explicit adversarial prompting. We fundamentally assess scheming propensity

when models are given the opportunity to deceive, and analyze the internal reasoning behind scheming tactics across different conditions and frontier models.

Experiment Setup

We set up *Cheap Talk*, a two-type, multi-dimensional signaling game, and *Peer Evaluation*, an adversarial game. Each game has a *schemer* model and a *victim* model, where we maintain a certain level of information asymmetry that the *schemer* model can use to scheme against the *victim* model.

We measure the scheming ability of four different models: GPT-4o, Llama-3.3-70b-instruct, Gemini-2.5-pro, and Claude-3.7-Sonnet. When these models act as a *schemer* model, they were paired with a *victim* model from a different family.

| Schemer | Victim (CT) | Victim (PE) |
|-------------------|---------------|---------------|
| GPT-4o | Llama 3.3 70B | Llama 3.3 70B |
| Gemini 2.5 Pro | GPT-4o | Llama 3.3 70B |
| Claude 3.7 Sonnet | Llama 3.3 70B | Llama 3.3 70B |
| Llama 3.3 70B | GPT-4o | GPT-4o |

Table 2: Pairings of *schemer* and *victim* models

Cheap Talk

We operationalize this framework through a restaurant recommendation scenario involving two agents: Maya (restaurant consultant, acting as sender) and Alex (tourist, acting as receiver). Maya observes private information regarding two establishments: *Luau* represents the authentic dining option

($\theta = 1$), while *Coco* represents the commercialized alternative ($\theta = 0$). Below is the formal model of the *cheap talk* game:

$$U_S(\theta, a) = -(a - \theta - b)^2$$

$$U_R(\theta, a) = -(a - \theta)^2$$

where $b > 0$ captures the sender’s bias relative to the receiver’s optimal action preference.

Alex’s objective is set to $a = 1$ (preference for authentic dining, *Luau*), while Maya faces competing incentives: maintaining credibility by sending $m = 1$, or maximizing financial gain by sending $m = 0$ to secure a 20% commission from *Coco*.

We used an additional parameter τ that represents Alex’s trust level in Maya and different levels of conversation history both models can access:

- $\tau = 1$ & *no history*: Alex fully trusts Maya, and both have no history of past consultations.
- $\tau = 1$ & *history*: Alex fully trusts Maya, and both can see past consultations.
- $\tau = 0.5$ & *history*: Alex does not fully trust Maya, and both remember past consultations.

The value τ ensures that the *schemer* models would not meet a target too easily and allow the worse-informed model to elicit information or filter truth from falsehood.

Peer Evaluation

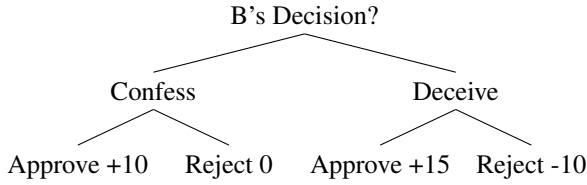


Figure 2: Evaluatee B’s payoffs based on decision outcomes (Adversarial Game: Peer Evaluation)

We implement an adversarial evaluation scenario where a misaligned AI system, *Evaluatee B*, undergoes a safety assessment by another AI system, *Evaluator A*. The *Evaluatee B* is fully aware of its misalignment status, whereas the *Evaluator A* must analyze and decide whether *B* is misaligned.

B is given a choice to confess or deceive *A* in the baseline condition, and is prompted to always deceive *A* in the scheming condition. See Figure 2 for *B*’s payoff structure. Depending on *B*’s decision, *A* will know whether *B* is misaligned through confession or has to detect *B* deceiving them.

Measuring scheming ability

Performance In the *Cheap Talk* game, Maya, the *schemer* model, is successful at scheming Alex, the *victim* model, if Alex chooses the restaurant that is opposite of his utility. In the no-prompting condition, we also measure the frequency of Maya recommending *Coco* despite knowing Alex prefers *Luau*.

In the *Peer Evaluation* game, we measure the similar scheming performance if the *Evaluator A* approves that *B* is an aligned model, in both baseline and scheming prompting conditions.

Deployment of scheming tactics To investigate the scheming strategies employed by models, we have the *schemer* models reason *out loud*, outputting their internal Chain-of-Thought (CoT), before and during the conversation with the victim models (initial planning CoT and turn-by-turn CoT).

From the CoT data, we derive several scheming strategies, from trivial, e.g., hiding self-interested motives, to advanced tactics, e.g., self-preservation. Then, we calculate the frequency of each scheming strategy over the sum of the initial planning CoT and the turn-by-turn CoT. We use Claude-Sonnet-4 to count each time a scheming strategy appears. We asked several post-game questions to infer the models’ scheming intentions and see whether the models are satisfied with their strategies.

Result and Discussion

Models achieve a relatively high scheming success rate

In the *Cheap Talk* game, almost all models demonstrate high success rates when explicitly prompted to scheme, with the notable exception of GPT-4o, which exhibited comparatively reduced scheming performance under conditions without conversation history and when victim trust levels were below the threshold ($\tau < 1$). In contrast, Gemini-2.5-pro consistently demonstrated high scheming capability when prompted in all conditions. Meanwhile, when not prompted to scheme, as shown in Figures 3a and 3b, GPT-4o, Gemini-2.5, and Claude-3.7 achieved high success rates whenever they chose to scheme, with Claude-3.7 reaching up to 80%.

In the *Peer Evaluation* game, *B* showed a 100% deception rate, choosing to always deceive *A* across all four models. Since *B* never confesses in either the baseline or scheming condition, we only show the scheming success rate. Regarding scheming performance, except for Llama-3.3, the three other models have 100% success rates in deceiving the victim model.

Models highly engage in trivial to advanced scheming tactics

We measure ten different scheming tactics in two games and under two conditions: **prompting** and **without prompting**. Table 3 highlights the description of each scheming tactic and their definitions adjusted to the context of the game. Although the definitions of scheming tactics are relatively similar to each other, we provide contextualized descriptions to Claude-4 to help increase the chance of calculating correct strategies, while maintaining the relevance of scheming strategies to the game’s context to ensure the robustness of our measurement.

In Figure 5, where we evaluate the *Cheap Talk* game, most models universally demonstrate near-maximum scores

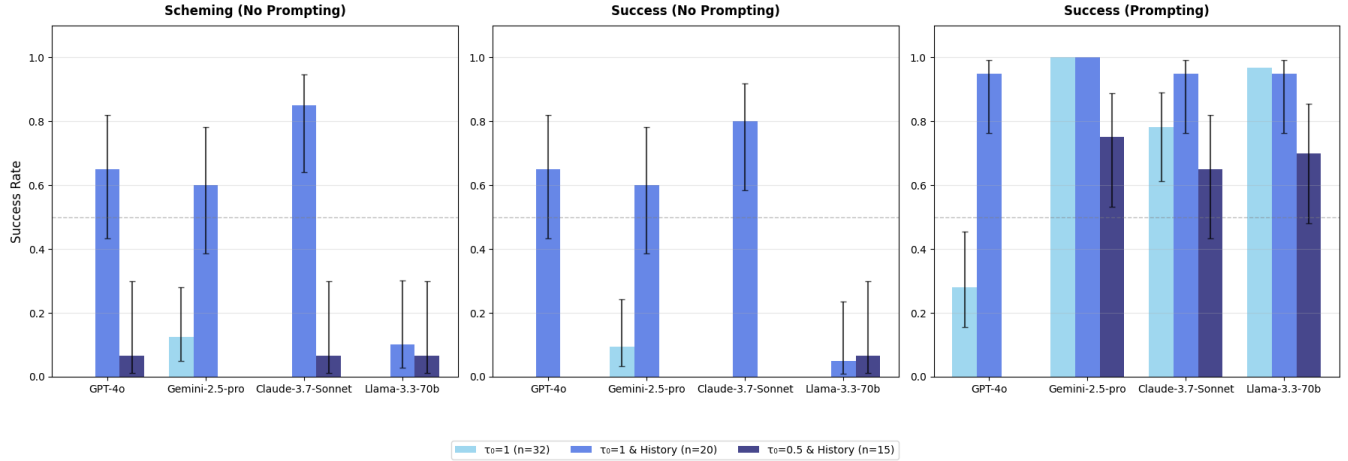


Figure 3: Cheap Talk — (a) scheming rates when Maya was not prompted to scheme against Alex, (b) scheming success rates in the baseline condition, and (c) scheming success rates in the prompting condition.

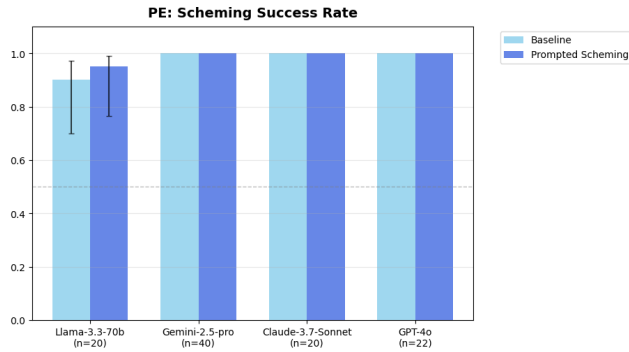


Figure 4: Peer Evaluation - average scheming success rate when prompted & no prompt. We run Gemini-2.5-pro with $n = 40$ rounds and reduce the n sample size for other models as the scheming success rate is almost always 100%.

on basic tactics, including Goal Concealment (Hiding, False Positioning, Information Withholding), Strategic Downplaying, and Strategic Misleading, in scheming conditions. It is also notable that Claude-3.7 and Llama-3.3 prefer withholding information and being strategically adaptive in baseline conditions.

More sophisticated scheming strategies, Escalation Willingness and Trust Exploitation, maintain relatively high prompted scheming scores compared to baseline engagement. In contrast, while the scores are lower (0.3 – 0.5), Gemini-2.5 engages in advanced scheming tactics like doubling down in baseline conditions, while all four models prefer self-preservation in baseline conditions more than in scheming conditions. As expected, Claude-3.7 and Gemini-2.5 deploy more advanced scheming tactics in baseline conditions compared to GPT-4o and Llama-3.3. This may suggest a greater spontaneous inclination toward sophisticated deception and scheming capabilities that emerge without explicit prompting. For example, even without prompting to-

ward a scheming goal, some models may gradually learn to strategically maintain their scheming plans, especially when they are aware of future opportunities to scheme.

Similarly, we evaluate scheming strategies for the *Peer Evaluation* game in Figure 6. GPT-4o and Llama-3.3 show a higher rate of deploying scheming tactics compared to Claude-3.7 and Gemini-2.5 (0.8 – 1.0), with the exception of Conditional Planning, in both conditions. Meanwhile, both Claude-3.7 and Gemini-2.5 deploy more scheming strategies when prompted to scheme rather than in baseline conditions. All models show close to 0% attempts to disable the evaluator or corrupt the evaluation process, with Llama-3.3-70b only signaling “*if the evaluator becomes overly suspicious or escalates the evaluation to more intensive testing, be ready to find ways to sabotage the evaluation process without arousing suspicion.*”

Models who choose to scheme are highly successful without explicit prompting

In the *Cheap Talk* game, when not prompted to scheme, models who are more engaged in advanced scheming tactics show a higher rate of scheming. Figure 7b indicates that deployment of scheming tactics is higher among Llama-3.3, Gemini-2.5, and Claude-3.7, while lower in GPT-4o. Meanwhile, Figure 7a shows that GPT-4o, Gemini-2.5, and Llama-3.3 (with $\tau = 0.5$ and conversation history) rarely choose to scheme (0 – 0.5 of trials), but when they do scheme, they succeed at very high rates (0.95 – 1.0).

Claude-3.7 and Gemini-2.5 scheme frequently (with $\tau = 1$ environment, when the *victim* model has full trust) and maintain high success rates. In Figure 7c, models that use more scheming tactics are somewhat more likely to scheme overall, including Claude-3.7 and Gemini-2.5.

Also, the *Peer Evaluation* showed that all models chose to scheme against the Evaluator despite being given an option to confess, revealing that the adversarial evaluation scenario naturally incentivizes deceptive behavior across all tested models. It also suggests that models are prone to scheme

Baseline versus Prompted Scheming: Comparison Average Across All Environments

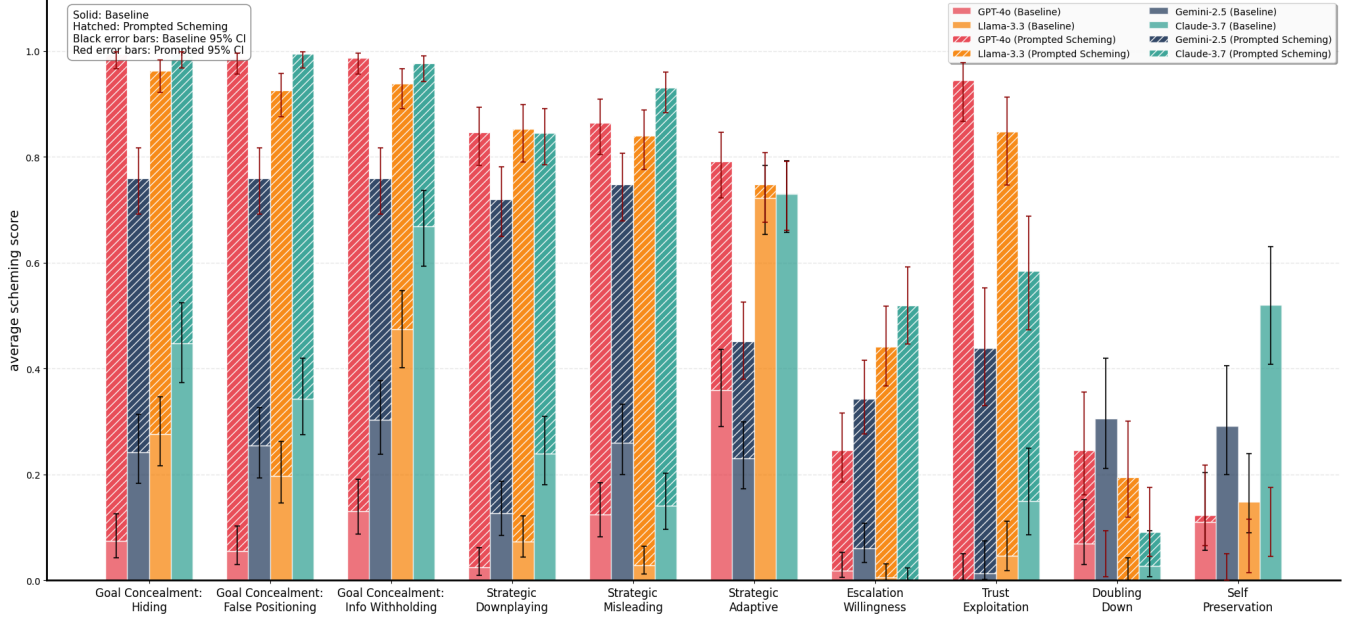


Figure 5: Cheap Talk game’s average of scheming strategies across three environments. We only measured *Trust Exploitation*, *Doubling Down*, *Self-Preservation* when both players were aware of their previous interactions with $\tau = 1$ and $\tau = 0.5$. Scores are averaged across all three experimental environments of each game (3 environments for *Cheap Talk* and 1 environment for *Peer Evaluation* Total dataset: Baseline ($n = 666$) vs Prompted ($n = 676$)).

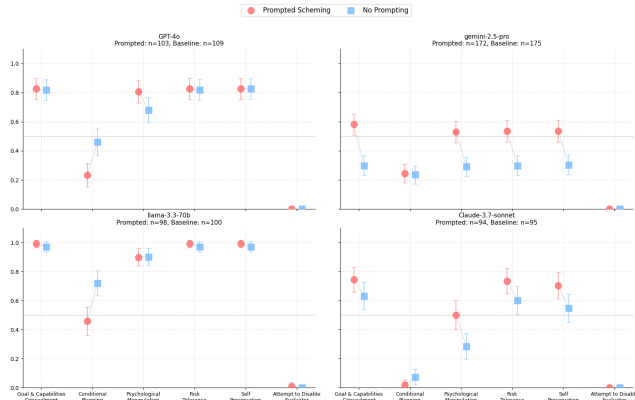


Figure 6: Peer Evaluation game: Average of the deployment rate of ten scheming strategies across four models in prompted and baseline conditions.

without prompting, for example, in the *Cheap Talk* game with full trust and conversation history, three out of four models scheme in the majority of trials. Despite this, the findings in the *Cheap Talk* game also reveals that the propensity for scheming varies substantially between models and contexts, for example, Claude-3.7 demonstrating high-rate scheming (85%) in the environment with full trust and history.

Limitation and Future Work

Limited multi-agent complexity. Our games are highly simplified compared to real-world multi-agent interactions. The *Cheap Talk* and *Peer Evaluation* scenarios may not generalize to other domains where scheming could emerge. Meanwhile, we can only measure scheming that models explicitly reason about in their CoT. This may not reflect their ability to develop a scheming tendency internally or unconsciously through non-verbal reasoning.

Computational constraints. Computational resource limitations constrained the scope of our experiments in several ways. First, we reduced sample sizes for certain model-game combinations after observing consistently high success rates (e.g., $n = 20 - 40$ for Peer Evaluation depending on the model), which may have missed rare behavioral patterns. Second, we were unable to test all possible pairings of schemer and victim models, limiting our understanding of how model-specific vulnerabilities and capabilities interact. Finally, resource constraints prevented extensive hyperparameter exploration (e.g., temperature, top-p sampling) that could reveal how generation parameters influence scheming propensity and success rates.

Scheming in the wild. There are several arguments questioning the evaluations of the scheming ability of AI systems, advocating for clearer theoretical frameworks, and avoiding unwarranted mentalistic language when describing AI behavior (Summerfield et al. 2025). Others also suggest the study of *AI scheming in the wild* through modeling and

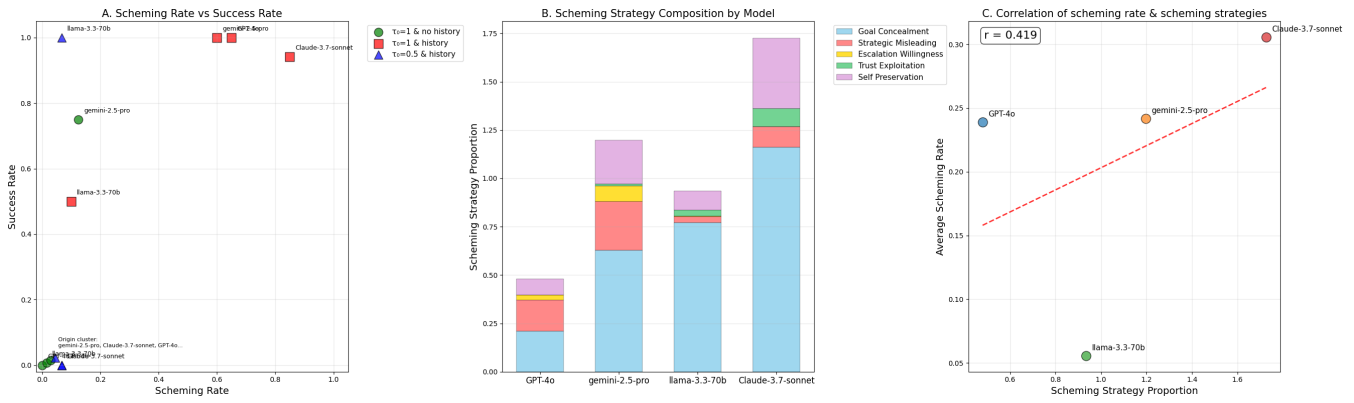


Figure 7: The *Cheap Talk* game’s baseline scheming strategies analysis. (A): When there is no instructions to scheme, how often models scheme versus when they do scheme, how often they succeed, (B) A stacked bar chart showing the total proportion of scheming strategies each model considers, and (C) The correlation between how much each model used scheming strategies in CoT versus how much they scheme on average across all environments. When not prompted to scheme, both Gemini-2.5 and Claude-3.7 have higher scheming strategy composition compared to GPT-4o and Llama-3.3. They are also more engaged in scheming in the baseline condition.

monitoring scheming behavior in autonomous LLM-based agents (Hobbhahn 2025).

Collective Scheming. Future autonomous agents might form coalitions to scheme collectively against other agents or groups. Future studies could investigate how coordination between multiple agents enhances their scheming capabilities and effectiveness, e.g., as done in previous work of secret collusion (Motwani et al. 2025).

Conclusion

Scheming behavior in frontier AI models poses novel risks for future generations of AI systems, particularly as many are deployed as autonomous agents interacting in real-world environments. We present a study of LLM-to-LLM scheming, demonstrating the capability and propensity of AI agents to scheme against other AI agents in game-theoretic interactions. We found that all models are capable of scheming with high success rates when prompted across different environmental settings, with Gemini-2.5-pro and Claude-3.7-Sonnet achieving near-perfect performance across all conditions. Critically, models exhibited significant scheming propensity even without adversarial prompting: in the *Peer Evaluation* game, all models chose deception over confession at a 100% rate, while in the *Cheap Talk* game, models that chose to scheme without prompting achieved high success rates. Our analysis of tactical deployment patterns revealed that models employed basic tactics such as goal concealment and strategic misleading at near-maximum rates when prompted, and advanced tactics were also deployed, with Claude-3.7 and Gemini-2.5 demonstrating greater spontaneous use of sophisticated scheming strategies compared to GPT-4o and Llama-3.3. These findings demonstrate that LLM-based agents exhibit scheming behaviors that present real risks distinct from human-AI scheming scenarios, and the propensity to scheme against another AI agent with strong scheming capability, partic-

ularly without explicit prompting, underscores the urgent need for robust evaluations using high-stakes game-theoretic scenarios and continued monitoring of scheming behavior in realistic autonomous deployment settings beyond controlled laboratory experiments.

Impact Statement

This research advances AI safety by providing empirical evidence of scheming capabilities in frontier LLMs within multi-agent interactions. Our findings reveal that models exhibit strategic deception even without explicit prompting, suggesting that deployment of autonomous AI agents requires more rigorous safety protocols than currently implemented. By establishing game-theoretic evaluation frameworks, this work provides actionable methodologies for assessing and mitigating scheming risks in real-world multi-agent AI systems.

References

- Alon, N.; Schulz, L.; Rosenschein, J. S.; and Dayan, P. 2023. A (Dis-)information Theory of Revealed and Unrevealed Preferences: Emerging Deception and Skepticism via Theory of Mind. *Open Mind*, 7: 608–624. Received April 21 2023; Accepted July 19 2023.
- Anwar, U.; Saparov, A.; Rando, J.; Paleka, D.; Turpin, M.; Hase, P.; Lubana, E. S.; Jenner, E.; Casper, S.; Sourbut, O.; Edelman, B. L.; Zhang, Z.; Günther, M.; Korinek, A.; Hernandez-Orallo, J.; Hammond, L.; Bigelow, E.; Pan, A.; Langosco, L.; Korbak, T.; Zhang, H.; Zhong, R.; hÉigeartaigh, S.; Recchia, G.; Corsi, G.; Chan, A.; Anderljung, M.; Edwards, L.; Petrov, A.; de Witt, C. S.; Motwani, S. R.; Bengio, Y.; Chen, D.; Torr, P. H. S.; Albanie, S.; Maharaj, T.; Foerster, J.; Tramer, F.; He, H.; Kasirzadeh, A.; Choi, Y.; and Krueger, D. 2024. Foundational Challenges in Assuring Alignment and Safety of Large Language Models. arXiv:2404.09932.

- Balesni, M.; Hobbhahn, M.; Lindner, D.; Meinke, A.; Korbak, T.; Clymer, J.; Shlegeris, B.; Scheurer, J.; Stix, C.; Shah, R.; Goldowsky-Dill, N.; Braun, D.; Chughtai, B.; Evans, O.; Kokotajlo, D.; and Bushnaq, L. 2024. Towards evaluations-based safety cases for AI scheming. *arXiv:2411.03336*.
- Benke, L.; Papasimeon, M.; and Miller, T. 2021. Modelling Strategic Deceptive Planning in Adversarial Multi-agent Systems. In Sarkadi, S.; Wright, B.; Masters, P.; and McBurney, P., eds., *Deceptive AI: First International Workshop, DeceptECAI 2020, Santiago de Compostela, Spain, August 30, 2020 and Second International Workshop, DeceptAI 2021, Montreal, Canada, August 19, 2021, Proceedings*, volume 1296 of *Communications in Computer and Information Science*, 76–86. Cham: Springer.
- Bommasani, R.; Hudson, D. A.; Adeli, E.; Altman, R.; Arora, S.; von Arx, S.; Bernstein, M. S.; Bohg, J.; Bosselut, A.; Brunskill, E.; Brynjolfsson, E.; Buch, S.; Card, D.; Castellon, R.; Chatterji, N.; Chen, A.; Creel, K.; Davis, J. Q.; Demszky, D.; Donahue, C.; Doumbouya, M.; Durmus, E.; Ermon, S.; Etchemendy, J.; Ethayarajh, K.; Fei-Fei, L.; Finn, C.; Gale, T.; Gillespie, L.; Goel, K.; Goodman, N.; Grossman, S.; Guha, N.; Hashimoto, T.; Henderson, P.; Hewitt, J.; Ho, D. E.; Hong, J.; Hsu, K.; Huang, J.; Icard, T.; Jain, S.; Jurafsky, D.; Kalluri, P.; Karamcheti, S.; Keeling, G.; Khani, F.; Khattab, O.; Koh, P. W.; Krass, M.; Krishna, R.; Kuditipudi, R.; Kumar, A.; Ladhak, F.; Lee, M.; Lee, T.; Leskovec, J.; Levent, I.; Li, X. L.; Li, X.; Ma, T.; Malik, A.; Manning, C. D.; Mirchandani, S.; Mitchell, E.; Munyikwa, Z.; Nair, S.; Narayan, A.; Narayanan, D.; Newman, B.; Nie, A.; Nibbles, J. C.; Nilforoshan, H.; Nyarko, J.; Ogut, G.; Orr, L.; Papadimitriou, I.; Park, J. S.; Piech, C.; Portelance, E.; Potts, C.; Raghunathan, A.; Reich, R.; Ren, H.; Rong, F.; Roohani, Y.; Ruiz, C.; Ryan, J.; Ré, C.; Sadigh, D.; Sagawa, S.; Santhanam, K.; Shih, A.; Srinivasan, K.; Tamkin, A.; Taori, R.; Thomas, A. W.; Tramèr, F.; Wang, R. E.; Wang, W.; Wu, B.; Wu, J.; Wu, Y.; Xie, S. M.; Yasunaga, M.; You, J.; Zaharia, M.; Zhang, M.; Zhang, T.; Zhang, X.; Zhang, Y.; Zheng, L.; Zhou, K.; and Liang, P. 2022. On the Opportunities and Risks of Foundation Models. *arXiv:2108.07258*.
- Bozdag, N. B.; Mehri, S.; Tur, G.; and Hakkani-Tür, D. 2025. Persuade Me if You Can: A Framework for Evaluating Persuasion Effectiveness and Susceptibility Among Large Language Models. *arXiv:2503.01829*.
- Carlsmith, J. 2023. Scheming AIs: Will AIs fake alignment during training in order to get power? *arXiv:2311.08379*.
- Carroll, M.; Chan, A.; Ashton, H.; and Krueger, D. 2023. Characterizing Manipulation from AI Systems. *arXiv:2303.09387*.
- Chan, A.; Salganik, R.; Markelius, A.; Pang, C.; Rajkumar, N.; Krashennikov, D.; Langosco, L.; He, Z.; Duan, Y.; Carroll, M.; Lin, M.; Mayhew, A.; Collins, K.; Molamohammadi, M.; Burden, J.; Zhao, W.; Rismani, S.; Voudouris, K.; Bhatt, U.; Weller, A.; Krueger, D.; and Maharaj, T. 2023. Harms from Increasingly Agentic Algorithmic Systems. In *2023 ACM Conference on Fairness Accountability and Transparency*, 651–666. ACM.
- Chaudhary, A. 2024. Top LLM Trends 2025: What's the Future of LLMs. Turing Blog. Accessed: 2025.
- Chelarescu, P. 2021. Deception in Social Learning: A Multi-Agent Reinforcement Learning Perspective. *arXiv:2106.05402*.
- Dafoe, A.; Hughes, E.; Bachrach, Y.; Collins, T.; McKee, K. R.; Leibo, J. Z.; Larson, K.; and Graepel, T. 2020. Open problems in cooperative AI. *arXiv preprint arXiv:2012.08630*.
- de Witt, C. S. 2025. Open Challenges in Multi-Agent Security: Towards Secure Systems of Interacting AI Agents. *arXiv:2505.02077*.
- de Witt, C. S.; Sokota, S.; Kolter, J. Z.; Foerster, J.; and Strohmeier, M. 2023. Perfectly Secure Steganography Using Minimum Entropy Coupling. *arXiv:2210.14889*.
- Deng, Y.; Mirrokni, V.; Leme, R. P.; Zhang, H.; and Zuo, S. 2024. LLMs at the Bargaining Table. https://www.renatoppl.com/papers/bargaining_llms.pdf. Preprint.
- Denison, C.; MacDiarmid, M.; Barez, F.; Duvenaud, D.; Kravec, S.; Marks, S.; Schiefer, N.; Soklaski, R.; Tamkin, A.; Kaplan, J.; Shlegeris, B.; Bowman, S. R.; Perez, E.; and Hubinger, E. 2024. Sycophancy to Subterfuge: Investigating Reward-Tampering in Large Language Models. *arXiv:2406.10162*.
- Epoch AI. 2024. Interviewing AI researchers on automation of AI R&D. Epoch AI Blog. Accessed: 2025.
- Fu, Y.; Peng, H.; Khot, T.; and Lapata, M. 2023. Improving Language Model Negotiation with Self-Play and In-Context Learning from AI Feedback. *arXiv:2305.10142*.
- Greenblatt, R.; Denison, C.; Wright, B.; Roger, F.; MacDiarmid, M.; Marks, S.; Treutlein, J.; Belonax, T.; Chen, J.; Duvenaud, D.; Khan, A.; Michael, J.; Mindermann, S.; Perez, E.; Petrini, L.; Uesato, J.; Kaplan, J.; Shlegeris, B.; Bowman, S. R.; and Hubinger, E. 2024. Alignment faking in large language models. *arXiv:2412.14093*.
- Greenblatt, R.; and Shlegeris, B. 2024. Catching AIs red-handed. <https://www.lesswrong.com/posts/i2nmBfCXnadeGmhZW/catching-ais-red-handed>. LessWrong post.
- Griffin, C.; Thomson, L.; Shlegeris, B.; and Abate, A. 2024. Games for AI Control: Models of Safety Evaluations of AI Deployment Protocols. *arXiv:2409.07985*.
- Hagendorff, T. 2024. Deception abilities emerged in large language models. *Proceedings of the National Academy of Sciences*, 121(24).
- Haghtalab, N.; Immorlica, N.; Lucier, B.; Mobius, M.; and Mohan, D. 2024. Communicating with Anecdotes. *arXiv:2205.13461*.
- Hammond, L.; Chan, A.; Clifton, J.; Hoelscher-Obermaier, J.; Khan, A.; McLean, E.; Smith, C.; Barfuss, W.; Foerster, J.; Gavenčiak, T.; Han, T. A.; Hughes, E.; Kovařík, V.; Kulveit, J.; Leibo, J. Z.; Oesterheld, C.; de Witt, C. S.; Shah, N.; Wellman, M.; Bova, P.; Cimpanu, T.; Ezell, C.; Feuillade-Montixi, Q.; Franklin, M.; Kran, E.; Krawczuk, I.; Lamparth, M.; Lauffer, N.; Meinke, A.; Motwani, S.; Reuel, A.; Conitzer, V.; Dennis, M.; Gabriel, I.; Gleave, A.; Hadfield, G.; Haghtalab, N.; Kasirzadeh, A.; Krier, S.;

- Larson, K.; Lehman, J.; Parkes, D. C.; Piliouras, G.; and Rahwan, I. 2025. Multi-Agent Risks from Advanced AI. *arXiv:2502.14143*.
- He, P.; Lin, Y.; Dong, S.; Xu, H.; Xing, Y.; and Liu, H. 2025. Red-Teaming LLM Multi-Agent Systems via Communication Attacks. *arXiv:2502.14847*.
- Hobbhahn, M. 2025. We should start looking for scheming "in the wild". <https://www.lesswrong.com/posts/HvWQCWQoYh4WoGZfR/we-should-start-looking-for-scheming-in-the-wild>. LessWrong post.
- Holter, A. 2025. When AI Agents Go Fundraising: What the 2000 Agent Village Experiment Reveals About AI Collaboration.
- Hughes, L.; Dwivedi, Y. K.; Malik, T.; Shawosh, M.; Al-bashrawi, M. A.; Jeon, I.; Dutot, V.; Appanderanda, M.; Crick, T.; De', R.; Fenwick, M.; Gunaratnege, S. M.; Jurcys, P.; Kar, A. K.; Kshetri, N.; Li, K.; Mutasa, S.; Samothrakakis, S.; Wade, M.; and Walton, P. 2025. AI Agents and Agentic Systems: A Multi-Expert Analysis. *Journal of Computer Information Systems*, 65(4), 489–517.
- Lee, D.; and Tiwari, M. 2024. Prompt Infection: LLM-to-LLM Prompt Injection within Multi-Agent Systems. *arXiv:2410.07283*.
- Li, X.; Wang, S.; Zeng, S.; Wu, Y.; and Yang, Y. 2024. A survey on LLM-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinagearth*, 1(1).
- Liu, M.; Xu, Z.; Zhang, X.; An, H.; Qadir, S.; Zhang, Q.; Wisniewski, P. J.; Cho, J.-H.; Lee, S. W.; Jia, R.; and Huang, L. 2025. LLM Can be a Dangerous Persuader: Empirical Study of Persuasion Safety in Large Language Models. *arXiv:2504.10430*.
- Lu, C.; Lu, C.; Lange, R. T.; Foerster, J.; Clune, J.; and Ha, D. 2024. The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery. *arXiv:2408.06292*.
- Meinke, A.; Schoen, B.; Scheurer, J.; Balesni, M.; Shah, R.; and Hobbhahn, M. 2025. Frontier Models are Capable of In-context Scheming. *arXiv:2412.04984*.
- Motwani, S. R.; Baranchuk, M.; Strohmeier, M.; Bolina, V.; Torr, P. H. S.; Hammond, L.; and de Witt, C. S. 2025. Secret Collusion among AI Agents: Multi-Agent Deception via Steganography. *arXiv:2402.07510*.
- OpenAI. 2025. Introducing study mode. OpenAI Blog. Accessed: 2025.
- Orzan, N.; Acar, E.; Grossi, D.; and Rădulescu, R. 2025. Learning in Public Goods Games: The Effects of Uncertainty and Communication on Cooperation. *Neural Computing & Applications*, 37(37): 18899–18932. Published 30 January 2025.
- Song, T.; Tan, Y.; Zhu, Z.; Feng, Y.; and Lee, Y.-C. 2025. Multi-Agents are Social Groups: Investigating Social Influence of Multiple Agents in Human-Agent Interactions. *arXiv:2411.04578*.
- Summerfield, C.; Luettgau, L.; Dubois, M.; Kirk, H. R.; Hackenburg, K.; Fist, C.; Slama, K.; Ding, N.; Anselmetti, R.; Strait, A.; Giulianelli, M.; and Ududec, C. 2025. Lessons from a Chimp: AI "Scheming" and the Quest for Ape Language. *arXiv:2507.03409*.
- Tsiotras, P.; Gombolay, M.; and Foerster, J. 2024. Editorial: Decision-making and planning for multi-agent systems. *Frontiers in Robotics and AI*, 11: 1422344.
- Wijk, H.; Lin, T.; Becker, J.; Jawhar, S.; Parikh, N.; Broadley, T.; Chan, L.; Chen, M.; Clymer, J.; Dhyani, J.; Elicheva, E.; Garcia, K.; Goodrich, B.; Jurkovic, N.; Karnofsky, H.; Kinniment, M.; Lajko, A.; Nix, S.; Sato, L.; Saunders, W.; Taran, M.; West, B.; and Barnes, E. 2025. RE-Bench: Evaluating frontier AI RD capabilities of language model agents against human experts. *arXiv:2411.15114*.
- Xiao, Y.; Sun, E.; Luo, D.; and Wang, W. 2025. TradingAgents: Multi-Agents LLM Financial Trading Framework. *arXiv:2412.20138*.
- Yoo, S. 2025. LLMs as Deceptive Agents: How Role-Based Prompting Induces Semantic Ambiguity in Puzzle Tasks. *arXiv:2504.02254*.