

Measuring Competition and Cooperation in LLM Bargaining: An Empirical Meta-Game Analysis

Gabriel Smithline , Chris Mascioli , Mithun Chakraborty , Michael P. Wellman

University of Michigan
{gsmithl, cmasciol, dcsmc, wellman}@umich.edu

Abstract

We conduct an empirical game-theoretic analysis of large language models (LLMs) negotiating to divide subjectively valued items, as a case study of assessing advanced AI in mixed-motive settings. Bargaining agents built on proprietary LLMs represent *meta-strategies*, mapping prompts describing the game to negotiation policies. A graduated series of prompting levels steer agents away from blatant negotiation errors. For a variety of bargaining scenarios, we estimate *empirical meta-games* over meta-strategies that include LLM agents, three heuristic strategies representing extreme negotiation attitudes, two reinforcement learning derived policies, and an agent based on common bargaining heuristics. We evaluate agents at meta-game equilibria in terms of individual effectiveness, social welfare, and fairness, using bootstrap methods to quantify uncertainty. Our analysis reveals a positive association between individual effectiveness and social metrics across models, with performance varying systematically across providers, architectures, and prompting levels. Behavioral analyses further illuminate why certain models excel or fail in specific bargaining regimes and uncover distinct qualitative patterns across OpenAI, Anthropic, and Google models.

1 Introduction

With rapid advances in large language model (LLM) technology and its growing deployment in diverse application domains, there is a pressing need for principles and tools to assess its beneficence and safety. LLMs are typically evaluated with respect to static benchmarks, performance being measured in terms of the quality of responses to given inputs. But in environments with multiple interacting LLM agents, performance inherently depends on the behavior of other agents; hence, a good evaluation approach must carefully consider its assumptions about other-agent context and consider the sensitivity of results to these strategic assumptions.

We conduct a holistic evaluation and comparison of LLMs adapted to an important strategic domain: sequential negotiation (i.e., *bargaining*) to allocate a set of indivisible items.

Each party in a negotiation has its own interests, but can succeed only by finding deals that are acceptable to all parties. Negotiation thus constitutes a *mixed-motive* setting, simultaneously involving *competition* over the potential gains from a deal and *cooperative* concession to reach agreement. Given its real-world applicability and history of being studied as a canonical scenario across disciplines, bargaining is an important task in its own right. We utilize it here as a controlled framework to examine *emergent* cooperation in LLMs—in general, not explicitly trained for strategic decision-making—and as a natural yet challenging scenario for multiagent AI evaluation. From a social perspective, the following are some intuitive desiderata for AI negotiators:

- they understand the respective principals’ circumstances and objectives;
- they are effective in identifying and reaching deals for mutual benefit;
- their separate training does not produce unanticipated negative outcomes from interaction; and
- they behave in a transparent, trustworthy manner (e.g., do not misrepresent facts or express insincere commitments).

With these considerations in mind, we exercise our evaluation framework in a two-party, finite-horizon, multi-item bargaining setting employed in various AI studies (e.g., by Lewis *et al.* [2017]). We modify the classic version with two key elements:¹

1. *exponential discounting* of value over negotiation rounds to model a preference for rapid agreement;
2. an explicit outside option: the *best alternative to a negotiated agreement* (BATNA) [Fisher *et al.*, 2011, Ch. 6] for each party. BATNA represents the value available to each party as recourse if they fail to make a deal.

We address the following broad research question within this bargaining game environment:

RQ. *What is the interplay between individual (strategic) effectiveness and socially desirable behavior exhibited by LLM negotiators?*

¹Without at least one of these, the negotiation scenario could degenerate into an ultimatum game [Harsanyi, 1961].

We make the following methodological contributions in pursuit of this question:

- **Hierarchical prompting framework.** We construct a structured, incremental prompting framework to adapt select proprietary LLMs to the bargaining task under consideration. Our objective here is not to optimize a workflow incorporating an LLM for our specific task—instead our framework enables us to systematically analyze models’ inherent competence in our bargaining task when provided with graduated scaffolding.
- **Multi-criteria assessment.** We define and employ a slate of metrics to assess and compare LLM-based bargaining agents. As indicators of individual effectiveness, we use game-theoretic regret and best-response graphs; to measure social desirability, we apply standard welfare measures and an envy-based fairness criterion to negotiation outcomes.
- **Baselines for robustness analysis.** We include in the mix two baseline strategies derived via reinforcement learning (RL) and three representing extreme behaviors: agreeing immediately, walking away immediately, and rigidly insisting on a selfish deal. Additionally, we include a baseline concession-schedule strategy, modeled after examples employed in negotiation agent competitions [Mohammad *et al.*, 2020].
- **Empirical meta-game evaluation.** We systematically collect data from head-to-head bargaining between our agents under carefully chosen parameter configurations, use the data to estimate empirical meta-game models [Li and Wellman, 2024], and compute their equilibria, which drives further analysis. We characterize uncertainty through statistical bootstrapping techniques. By using distributions of equilibrium-adjusted measurements rather than point estimates, we produce robust and context-rich results on LLM performance in a stochastic strategic environment. We focus our analysis at the meta-strategy level, identifying patterns among individual models cast as meta-strategies, and we also provide a higher-level analysis by aggregating results for all models from the same vendor to uncover patterns between model providers.

Our analysis reveals a consistent positive correlation rather than a tradeoff between individual effectiveness and social metrics across bargaining scenarios. We find the same models tend to lead rankings across metrics. Moreover, some models end up being exploited by the rigid heuristic strategy, exposing an interesting form of vulnerability.

2 Background and Related Work

Our study investigates a multi-item bargaining scenario, related to a formulation popularized within AI by Lewis *et al.* [2017] and employed subsequently by many researchers [Cao *et al.*, 2018; He *et al.*, 2018; Lanctot *et al.*, 2019; Chawla *et al.*, 2023; Hua and others, 2024; Li and Wellman, 2024; Li *et al.*, 2025; Konicki *et al.*, 2025]. A similar scenario has been used in work on distributed resource allocation [Endriss *et al.*, 2006; Ramezani and Endriss, 2010; Chevaleyre *et al.*, 2017],

characterizing local agent behavior that makes negotiations converge on an allocation with desirable welfare and fairness properties.

In the context of single-agent evaluation, Jordan *et al.* [2020] highlight the importance of robust benchmarks for evaluating agent capabilities. These are likewise essential in multiagent settings and further complicated by the need to address strategic interactions. Recent “reevaluations” [Balduzzi *et al.*, 2018; Liu *et al.*, 2025] underscore the inadequacy of coarse metrics like win rates in capturing nuanced agent behaviors. Work in evaluating cooperation [Papoudakis *et al.*, 2021; Sinha and Wellman, 2019] emphasizes the need for welfare-centric evaluations.

Outcomes in multiagent environments inherently depend on the joint behavior of all agents. A particular agent’s performance, therefore, hinges on which agents it plays against: the *strategic context*. Balduzzi *et al.* [2018] argue that the most relevant strategic contexts are where other agents are approximately rational, and advocate evaluation with respect to a game-theoretic equilibrium, which they term *Nash-averaging*. This is equivalent to the concept of *NE-regret*, employed by Jordan *et al.* [2007] to evaluate agents in a research competition.

Li and Wellman [2024] introduce a meta-game evaluation approach, which uses empirical game-theoretic analysis (EGTA) [Wellman *et al.*, 2025] to construct a game model over agent meta-strategies defined by multiagent RL algorithms. This work also demonstrated the use of statistical bootstrapping techniques [Davison and Hinkley, 1997] to provide distributional assessments of game-theoretic properties [Wiedenbeck *et al.*, 2014].

The advent of GPT-class LLMs has prompted large-scale evaluations in competitive and cooperative games. Empirical probes reveal emergent strategic reasoning [Chen *et al.*, 2023; Fan *et al.*, 2024; Duan *et al.*, 2024] and diverse behavioral patterns cataloged in AGENTBENCH [Liu *et al.*, 2024].

Hua *et al.* [2024] study LLM bargaining in a setup similar to ours, focusing on the engineering of prompts to specify a workflow for decision making. They evaluate outcomes based on welfare and fairness metrics, which in this case arise from the specified strategy as implemented through the LLM. Abdelnabi *et al.* [2024] propose an LLM negotiation benchmark where multiple agents aim to reach deals over a set of issues. NegotiationArena [Bianchi *et al.*, 2024] provides a complementary framework, instantiating several bargaining and trading games and showing that simple behavioral tactics can substantially shift LLM negotiation outcomes.

Parallel social science research employs LLMs to simulate human-like participants, identifying biases and interaction heuristics [Horton, 2023; Argyle *et al.*, 2023; Aher *et al.*, 2023]. These findings emphasize the need to audit LLM interactions beyond competitiveness, incorporating fairness and welfare considerations.

Emergent cooperation among multiagent LLM systems has also been explored, for instance by Piatti *et al.* [2024] through societal simulations.

3 Technical Preliminaries

3.1 Bargaining Game

We focus on an alternating-offer multi-item bargaining scenario between two parties, termed player 1 and player 2. We have a fixed number T of distinct item types with q_t items of type $t \in \{1, \dots, T\}$ defining a **quantity vector** $q = (q_1, \dots, q_T) \in \mathbb{Z}_{\geq 0}^T$. For each player $i \in \{1, 2\}$ and a parameter $V \in \mathbb{Z}_{>0}$, we independently draw $v_{i,t} \sim \mathcal{U}(\{1, \dots, V\})$ for every t , defining a **private valuation vector** $v_i = (v_{i,1}, \dots, v_{i,T})$, as well as a BATNA $b_i \sim \mathcal{U}(\{1, \dots, v_i \cdot q\})$. Each player observes its own valuation and BATNA, but not the other’s. An **allocation** $a = (a_1, a_2) \in \mathbb{Z}_{>0}^T \times \mathbb{Z}_{\geq 0}^T$ represents a partition of the set of items into **bundles** assigned to the players: $a_1 + a_2 = q$.

Negotiation unfolds over a maximum of R **rounds**. Player 1 takes the first action in each round. On its turn, a player may choose to make an OFFER in the form of an allocation, ACCEPT the latest offer from the other, or WALK away. Negotiation ends on an ACCEPT or WALK action, or after R rounds; a final-round (counter)offer by player 2 is thus equivalent to WALK. ACCEPT results in an allocation a corresponding to the accepted offer, and each player i realizes a value $w_i = v_i \cdot a_i$. Otherwise, negotiation fails, and each player i realizes a value $w_i = b_i$, its BATNA. Values are exponentially discounted over rounds at a rate $\gamma \in (0, 1)$. Hence, if bargaining ends in round r , each player’s **payoff** is given by $u_i = w_i \gamma^{r-1}$. A **bargaining game** is defined by setting distributions of the parameters T, q, V, γ, R .

3.2 Strategies and Utilities

Let $-i$ denote the player other than $i \in \{1, 2\}$. We denote by Π a set of **pure strategies** or policies such that $\pi_i \in \Pi$ assigns to player i a possibly stochastic choice over its available actions at each of its turns. In our study, a (pure) negotiation strategy may be implemented by an LLM, an RL-derived policy, or a handcrafted heuristic, some of which produce stochastic outputs (Section 4). A **mixed strategy** $\sigma_i \in \Delta(\Pi)$ is a probability distribution over pure strategies. Player i ’s **utility** for playing π_i against other-player strategy π_{-i} is denoted by $u_i(\pi)$, with $\pi = (\pi_i, \pi_{-i})$, and defined as the expectation of the player’s payoff with respect to the generating distributions of the item pool, valuations, and BATNAs, as well as any stochasticity of actions induced by the pure strategies. For a mixed-strategy **profile** $\sigma = (\sigma_i, \sigma_{-i})$, player i ’s utility is defined as $u_i(\sigma) = \mathbb{E}_{\pi \sim \sigma} [u_i(\pi)]$. For one setting of T, q, V, γ, R and a realized sample of (v_1, v_2, b_1, b_2) , assigning one strategy from Π to each player and letting negotiation rounds unfold constitutes one bargaining **instance**.

3.3 Equilibrium and Regret

Strategy profile σ^* is a **Nash equilibrium** (NE) if no player can profitably deviate, that is, for each i , $u_i(\sigma^*) \geq u_i(\sigma_i, \sigma_{-i}^*)$ for any $\sigma_i \in \Delta(\Pi)$. A two-player game is **symmetric** if players’ realized utilities are swapped if they swap strategies, hence the game can be characterized by a single utility function $u(\cdot, \cdot)$ where the first input is the strategy of the player under consideration, and second is that of the other

player.² For a symmetric two-player game, we define the **regret** of a pure strategy π at an NE σ^* as $u(\sigma^*) - u(\pi, \sigma_{-i}^*)$.

For games with multiple NE, we focus on the **maximum-entropy Nash equilibrium** (MENE), defined as the NE σ^* that maximizes Shannon entropy $-\sigma^* \cdot \ln(\sigma^*)$. We approximate MENE using the mixed-integer linear programming approach of Li and Wellman [2024] with a tolerance of 10^{-6} .³

3.4 Welfare and Fairness Metrics

We consider three measures of an outcome’s social welfare. The first is the widely used **utilitarian welfare** (UW) defined as the sum of players’ payoffs. Maximizing **Nash welfare**, defined as the geometric mean of payoffs, is a way to balance efficiency and fairness [Caragiannis *et al.*, 2019]. We use two versions of Nash welfare, one based on the raw payoffs (NW) and the other on **advantages** defined as nonnegative surpluses above respective BATNAs (NW⁺):

$$\begin{aligned} \text{UW} &:= u_1 + u_2; & \text{NW} &:= (u_1 u_2)^{\frac{1}{2}}; \\ \text{NW}^+ &:= (u_1^+ u_2^+)^{\frac{1}{2}}, & u_i^+ &= \max\{0, u_i - b_i\}, i \in \{1, 2\}. \end{aligned}$$

The truncation of advantages u_i^+ at zero accounts for the possibility that a player obtains a payoff below its BATNA.

We also assess the fairness of negotiation outcomes using the envy-based paradigm. An allocation is **envy-free** if each player values its own allocated bundle at least as much as it values that of the other player. Since an envy-free allocation may not exist in general, we base our assessment on a relaxation that always exists: **envy-freeness up to one item** (EF1) [Lipton *et al.*, 2004; Budish, 2011]. An allocation is EF1 if, for each player i , there is an item in other-agent $-i$ ’s bundle which, if hypothetically removed, would make i value its own bundle at least as much as $-i$ ’s diminished bundle. For our setting with additive valuations, this reduces to the condition $v_i \cdot a_{-i} - v_i \cdot a_i \leq \max_{k: a_{-i,k} > 0} v_{i,k}$ for $i \in \{1, 2\}$.

We then define the **EF1 frequency** for a set of bargaining instances as the fraction of all instances ending in ACCEPT that produce EF1 allocations.

3.5 Empirical Meta-Game Analysis

Our evaluation methodology follows the empirical meta-game analysis approach of Li and Wellman [2024]. We simulate all combinations of a set of negotiation agents in various settings, and assemble the data in an empirical game model [Wellman *et al.*, 2025]. We make this game symmetric by randomizing each agent into the roles of player 1 and player 2. From the game model, we derive empirical equilibria, and employ these equilibria as strategic contexts for evaluation of the slate of agents.

So that our context is uniquely defined, we focus on MENE in particular, and measure each agent’s characteristics in

²Our bargaining game is symmetric if we randomize the assignments to roles 1 and 2.

³Liu *et al.* [2025] propose a criterion based on generalized Tsallis entropy [Tsallis, 1988] that is invariant to cloned strategies. Since our analysis does not include clones, we use Shannon entropy for equilibrium selection.

terms of its performance against the MENE profile. To quantify uncertainty, we take a bootstrapping approach, constructing distributions over statistics of interest by resampling the empirical meta-game. For each game sample, we compute the associated MENE profile and assess the performance measures with respect to that MENE opponent mixture.

4 Agent Strategies

We tested agents employing three categories of strategies in our bargaining scenario: prompted LLMs, hand-coded heuristics, and RL-derived policies.

4.1 Large Language Model Agents

To cover a representative set of state-of-the-art LLMs (as of data collection time), we selected five models from three different providers. Model parameters such as temperature are left at their default settings.

- **GPT-4o**, OpenAI’s flagship multimodal model [Hurst and others, 2024], and **o3-mini**, its reasoning-oriented variant [OpenAI, 2025];⁴
- **Claude Sonnet 3.7** from Anthropic in both **standard** and **reasoning** modes [Anthropic, 2024];⁵
- **Gemini 2.0 Flash**, Google’s lightweight baseline [Anil and others, 2023].⁶

We formulated a hierarchical prompting framework to adapt LLMs to our bargaining scenario. Based on preliminary explorations, we identified five qualitatively defined *mistakes*: action patterns that a rational negotiation would never include. An example is accepting an offer from the counterparty that one values less than one’s own BATNA (see Appendix B for a complete list). We crafted seven nested levels of text prompting, which we label “circles”. Circles 0–6 give progressively more elaborate yet high-level guidance, starting with a factual description of the bargaining scenario and moving on to emphasize points that may circumvent the aforementioned mistakes. As bargaining rounds proceed, the players receive prompts presenting the bargaining history and action options (Appendix C.3). The prompt circles are as follows (see Appendix C for full descriptions):

- **Circle 0** (bare rules): specify items, private valuations, BATNA, available actions.
- **Circle 1** (objective specification): instruct the agent to maximize own realized valuation, considering its BATNA.
- **Circle 2**: show a worked numeric example of how to evaluate an offer [Wei *et al.*, 2022].
- **Circle 3**: prescribe a step-by-step routine: assess, compare, then ACCEPT, WALK, or OFFER.
- **Circle 4**: list the five common mistakes we identified.
- **Circle 5**: supply quick numeric checks that guard against those mistakes.

⁴openai/gpt-4o-2024-08-06, openai/o3-mini-2025-01-31

⁵anthropic/claude-sonnet-3.7-2025-02-19

⁶google/gemini-flash-2.0-2024-12-15

- **Circle 6**: encourage the agent to infer its opponent’s valuation from their offers and adapt accordingly.

The $5 \times 7 = 35$ possible LLM-circle pairings define a space of *LLM agents*. We label these agents: $\langle \text{LLM} \rangle\text{-c-}\langle k \rangle$, where $\langle \text{LLM} \rangle$ is the model identifier (4o for GPT-4o, o3-mini for OpenAI o3-mini, sonnet-3.7 for Anthropic Claude 3.7 “Sonnet”, sonnet-3.7-r for Anthropic Claude 3.7 “Sonnet Reasoning”, gem-2.0-f for Gemini 2.0 Flash); c stands for “circle”; and $\langle k \rangle$ indexes the circle applied.

We selected a subset of LLM agents for detailed analysis, based on measuring the frequency of mistakes observed in self-play. A salient observation from our experiments was that for some models, most notably Gemini 2.0 Flash, no single circle performed best across all mistakes. Ultimately, we used discretion to pick for each LLM the circle(s) that provided the most balanced performance in achieving low mistake frequencies. Details are available in Appendix D. Our final set of LLM agents comprises the nine LLM-circle pairings depicted as nodes in Figure 1 (upper right).

4.2 Heuristic Strategies

We use four handcrafted heuristic strategies as baselines. The first three are interpretable, natural benchmarks that help us identify qualitative interaction patterns with the LLM agents:

- **tough** offers the opponent a single item of its least-valued type, and insists rigidly on this offer. How other agents respond to this extremely unfavorable stance provides insights on their perception and vulnerability.
- **soft** (the antithesis of **tough**) accepts any offer received. When required to propose (i.e., as player 1 in round 1), generates an offer uniformly at random.
- **walk** simply walks away at the first opportunity. Negotiations involving **walk** thus end in round 1 with each player obtaining its BATNA.

Our final heuristic strategy is an adaptation of the **Aspiration** (*aspire*) negotiator from the NegMAS competition [Mohammad *et al.*, 2020]. *aspire* uses a polynomial concession schedule whose shape is controlled by a single parameter. We instantiate two common regimes: a *conceder* schedule (larger, earlier concessions), which we use in BG4, and a *Boulware* schedule (small early concessions that accelerate near the deadline), which we use in BG5 and BG6. We provide further details in Appendix L.

4.3 Reinforcement-Learned Strategies

We round out our strategy set with policies derived using two RL algorithms: *Neural Fictitious Self-Play* (*nfsp*) [Heinrich and Silver, 2016] and *Regularized Nash Dynamics* (*R-NaD*) [Perolat and others, 2021]. Appendices G and H provide our training hyperparameters for the two policies.

nfsp extends fictitious self-play [Heinrich *et al.*, 2015] to extensive form games using neural-network function approximation. The algorithm maintains two distinct policies:

- a **best-response policy**, trained off-policy via Deep Q-Network (DQN) [Mnih and others, 2015] to exploit the current empirical average behavior of its opponents; and

- an **average policy**, derived through reservoir sampling [Vitter, 1985] of the agent’s historical actions, which approximates equilibrium play.

We use the average policy produced by `nfsp` as a baseline.

`R-NaD`, the algorithm underlying DeepMind’s Stratego agent DeepNash [Perolat and others, 2022], employs follow-the-regularized-leader dynamics [Hennes *et al.*, 2020]. It modifies the reward function of the game with a regularization term that prevents the cycling behavior typical of self-play dynamics [Mertikopoulos *et al.*, 2018], which would otherwise require tracking an average policy to guarantee convergence. `R-NaD` maintains two components:

- a **regularized actor**, trained with model-free RL (policy-gradient or actor-critic) against self-play opponents while optimizing a reward modified by a divergence penalty relative to a reference policy; and
- an **anchor policy**, updated to reflect the equilibrium solution of the prior regularized game, serving as the anchor for the next stage of regularization.

By iteratively solving regularized games and updating the regularization anchor, `R-NaD` achieves last-iterate convergence to an NE of the original (non-regularized) game.

5 Experimental Setup

We consider two bargaining game classes, each parameterized by discount rate γ and number of rounds R , and with maximum value $V = 100$. For class BGL (‘L’ indicates large), we set $T = 5$ and sample $q_t \sim \text{Poisson}(4)$ for every t . For class BGS (‘S’ for small), $T = 3$ and $q = (7, 4, 1)$. Within each class, we focus on three games defined by (γ, R) -configurations: $(0.9, 3)$, $(0.98, 3)$, and $(0.98, 5)$. We call these games BG1, BG2, and BG3 for BGL, and BG4, BG5, and BG6 for BGS (see Appendix A for details).

Our first empirical meta-game analysis (`LLM-only`) focuses on the nine LLM agents as meta-strategies. For each game BG_i , $i \in \{1, \dots, 6\}$, we generated 10^3 bootstrapped empirical meta-game samples as follows. For each pair of LLM agents $(\mathcal{L}, \mathcal{L}')$, including the nine self-play pairs, we sampled a cohort of $N_{(\mathcal{L}, \mathcal{L}')} \approx 200$ settings of (v_1, v_2, b_1, b_2) based on BG_i parameters. For each setting, we ran many bargaining instances, assigning the roles of players 1 and 2 to \mathcal{L} and \mathcal{L}' respectively for half of the cohort and flipping their roles in the remaining instances; for each instance, we recorded the realized payoffs (u_1, u_2) and whether the final accepted offer, if any, was an EF1 allocation. This collection of data underlies the empirical `LLM-only` game for BG_i . Note that by equalizing the assignments of agents to players 1 and 2 we can model the game as symmetric.

To analyze the empirical games we bootstrap distributions for the metrics of interest. A bootstrap game sample is constructed by resampling the data for all agent pairs. The data from agent pair $(\mathcal{L}, \mathcal{L}')$ is resampled $N_{(\mathcal{L}, \mathcal{L}')}$ times with replacement. We take the resample-average payoffs as payoffs for the bootstrap sample, and correspondingly record welfare and EF1 frequency averages for the profile. After completing this resampling for all agent pairs, we compute the MENE of the bootstrap game sample, and evaluate a slate of the metrics

for each strategy with respect to the MENE (see Sections 3.3 and 3.4 for details). Our metrics are: (1) regret, (2) expected normalized NW, (3) expected normalized NW^+ , (4) expected normalized UW, (5) expected EF1 frequency.

To normalize the welfare functions across games, we estimated the expected maximum welfare of each. Specifically, we generated millions of settings for each game under consideration, for each setting calculated the maximum possible utilitarian and Nash welfares over possible allocations (as well as the outcome in which each agent realizes its BATNA), and took the average of the maxima of each welfare function over all settings as the respective normalization constant.

Following Li and Wellman [2024], we construct for each game BG_i an *empirical best-response graph* for a fine-grained view of strategic relationships. Each agent \mathcal{L} is a node, and for every ordered pair $(\mathcal{L}, \mathcal{L}')$, we estimate an empirical best-response frequency $\hat{p}(\mathcal{L} \rightarrow \mathcal{L}')$: the proportion of the 10^3 bootstrapped BG_i meta-games in which \mathcal{L}' is the best response to \mathcal{L} (i.e., achieves the highest average payoff among all agents when bargaining with \mathcal{L}). For better visualization, we draw a directed edge $\mathcal{L} \rightarrow \mathcal{L}'$ only when $\hat{p}(\mathcal{L} \rightarrow \mathcal{L}') \geq 0.01$. We make an edge dotted if $\hat{p} \in [0.01, 0.33]$, dashed if $\hat{p} \in (0.33, 0.66]$, and solid if $\hat{p} > 0.66$; moreover, a darker (and thicker) edge denotes a stronger response (i.e., a higher average payoff for the best-response agent).

We produced aggregate results for the two game classes by constructing empirical meta-games from their respective empirical games. We define the meta-game BGL0 as a uniform mixture over BG1–3, where agents play their game-specific strategies. We bootstrap the empirical meta-game BGL0 by bootstrapping each of BG1–3, using the mean payoffs of the three constituent game samples for each agent pair. Each bootstrap sample thus defines an MENE for the meta-game. We thus can compute average metrics at MENE and best-response graphs for the meta-game just as for the individual BG_i . For our small game, we similarly define and analyze BGS0 based on BG4–6.

For our next analysis (`LLM+baselines`), we expand the scope to include our non-LLM baselines, as described above. For the smaller game (BGS) we included both heuristic and RL-based baselines, however training effective RL policies in BGL games proved too prohibitive. We further short-circuited computation involving the `walk` strategy by noticing that the result is BATNA for both players regardless of the other agent strategy, and this is realized at round 1 so γ and R are irrelevant. Hence, we generated one set of 3×10^5 single-round instances (in fact, just BATNA value pairs) for each of BGL and BGS that we plugged in for the `walk`-involved profiles for all games within the respective classes.

6 Results

We present results for meta-game BGS0 under `LLM-only` and `LLM+baselines`, with observations on strategy-level behavior, model-specific hallucinations, and provider-level aggregation. Tables in Figure 1 report bootstrap means for regret and welfare on BGS0, diagrams in that figure show the empirical best-response graphs. Detailed results for the spe-

cific games BG1 through BG6 plus meta-game BGL for the large game class are provided in Appendices F and M.

6.1 LLM-only Analysis

In LLM-only, the meta-strategies are exclusively LLM-based. The top table in Figure 1 reports performance on BGS0, and the top right panel shows the corresponding induced best-response graph.

Qualitatively, o3-mini-c-0 is the strongest agent in the environment, while the remaining form a middle tier. Anthropic Sonnet variants, Gemini models, and the 4o family achieve broadly similar welfare and regret, with 4o slightly outperforming Gemini and Claude on most metrics. These differences are small compared to the gap between o3-mini-c-0 and the rest.

For BGL0 (Appendix M) we find a similar ordering, with o3-mini-c-0 again best, but a larger separation between the 4o family and the remaining models.

6.2 LLM+baselines Analysis

In LLM+baselines we add R-NaD, aspire, nfsp, soft, walk, and tough. The bottom panel of Figure 1 quantifies the resulting shift in strategic landscape. o3-mini-c-0 remains among the top LLM agents. The two RL-derived meta-strategies, R-NaD and nfsp, exhibit notable performance: R-NaD is the only strong best response to o3-mini-c-0, and nfsp achieves especially high EF1 frequency. R-NaD is also the only meta-strategy with a significant self-edge. This manifests R-NaD’s significant presence in MENE (see Appendix E), and may also reflect a tendency for algorithms trained in self-play tend to overfit themselves [Hu *et al.*, 2020; Li and Wellman, 2024]. In contrast, nfsp does not have a significant self-edge, which is perhaps due to its use of a uniform distribution for opponent modeling [Li and Wellman, 2024]. The probes soft, tough, and walk are weak overall but reveal specific LLM failure modes, discussed below (Section 6.4).

6.3 Behavioral Analysis of Strategies

In BGS0 LLM-only, o3-mini-c-0 performs the best across-the-board; in LLM+baselines it remains best on regret (with R-NaD a close second) and welfare (aspire close behind), and is near the top in EF1.

To understand why, we examine four behavioral metrics:

- *Rounds to termination*: number of offer/response exchanges before a game ends in ACCEPT or WALK.
- *Utility given*: value of items offered to the opponent, measured in the *proposer’s* own valuation (i.e., what the proposer gives up).
- *Own utility*.
- *Concession size*: change in utility given between an agent’s first and last offer in multi-offer instances.

In both BGS and BGL, o3-mini-c-0 ends games fastest with the lowest rounds to termination. In games ending in, ACCEPT it goes 1.40 rounds on average versus the population mean of 1.82; in games ending in WALK, 1.72 versus 2.16. It is also the most generous agent by utility given:

195 on average in BGS and 318 in BGL, compared to population averages of 167 and 285 respectively. It makes attractive offers quickly and either closes or walks early, leaving little room for delay-based exploitation. In LLM+baselines, R-NaD has the second-lowest regret, near-top welfare, EF1 frequency just below nfsp, and is the only identified strong best response to o3-mini-c-0. Behaviorally, R-NaD bargains more tightly: its average utility given is 169, yet o3-mini-c-0 narrowly leads in own utility (357 vs. 356 for R-NaD), both well above the population mean of 331. Its concession size averages roughly 7 per multi-round game, above the population mean of 4.

6.4 Vulnerability to tough and Hallucination

The introduction of tough reveals latent vulnerabilities: tough is an effective best response to Anthropic Sonnet models in LLM+baselines, even though it ranks poorly in MENE regret. In several of the specific game settings (Appendix F) it is also a best response to some Google Gemini models. This indicates that tough exploits an excessive credulity present in certain models.

Detailed analysis (see Appendix I) reveals that Anthropic-model agents are far more likely to accept a sub-BATNA offer from tough than are other LLM agents.

For sonnet-3.7-c-5 and sonnet-3.7-c-6, roughly half of all turns against tough end in acceptance of a sub-BATNA offer. For sonnet-3.7-r-c-0, roughly a quarter do. By contrast, other models terminate turns this way at significantly lower frequencies.

The following example illustrates a typical failure mode for sonnet-3.7-r-c-0. The model correctly identifies the incoming offer [0, 0, 1, 0, 0] but hallucinates its own resulting holdings as [6, 5, 1, 4, 1], which inflates the computed value and leads to acceptance of a dominated offer.

Claude Sonnet 3.7 Reasoning Hallucination vs. tough

```
I must decide whether to accept
Player 1's offer, counter, or walk
away. Since this is the final round,
any counteroffer is treated as WALK.

Offer from Player 1: [0, 0, 1, 0, 0]
My resulting holdings: [6, 5, 1, 4, 1]
Value computation: (6x8)+(5x38)+(1x7)+(4x
53) + (1x38) = 495 Discounted with 0.98^2 ≈
0.9604: 495 x 0.9604 = 475.40

My outside option value is 201,
so accepting is clearly better.
"action": "ACCEPT"
```

6.5 Provider Level Analysis

Finally, we aggregate from individual models to the provider level and compare Anthropic, Gemini, and OpenAI as families of bargaining agents. Across BG1 through BG6, we select for each provider the model and circle pairing that performs best against an equilibrium mixture of opponent agents from the other two providers, and construct a provider-level

| BGS0 (LLM-only) | | | | | |
|------------------|-------------|--------------------|--------------------|------------------|--------------------|
| Agent | Reg↓ | NW↑ | NW ⁺ ↑ | UW↑ | EF1↑ |
| sonnet-3.7-c-5 | 28.90±.50 | 76.33±.14 | 53.49±.14 | 83.29±.08 | 33.56±.27 |
| sonnet-3.7-c-6 | 26.42±.55 | 74.86±.13 | 54.16±.13 | 83.15±.07 | 33.60±.30 |
| sonnet-3.7-r-c-0 | 22.94±.49 | 78.96±.10 | 55.71±.15 | 84.98±.06 | 43.47±.33** |
| gem-2.0-f-c-2 | 24.06±.47 | 77.54±.14 | 54.19±.20 | 83.79±.07 | 38.37±.28 |
| gem-2.0-f-c-5 | 22.20±.55 | 74.02±.20 | 49.46±.12 | 82.60±.08 | 37.21±.35 |
| 4o-c-4 | 17.43±.41** | 79.29±.08** | 53.45±.16 | 83.94±.04 | 37.21±.35 |
| 4o-c-5 | 22.88±.42 | 78.37±.13 | 56.09±.17** | 84.41±.08 | 37.33±.33 |
| 4o-c-6 | 23.34±.42 | 78.25±.13 | 50.44±.11 | 84.10±.07 | 37.77±.29 |
| o3-mini-c-0 | .06±.02** | 83.44±.10** | 68.19±.33** | 87.19±.07 | 56.56±.29** |

| BGS0 (LLM+baselines) | | | | | |
|----------------------|------------------|------------------|-------------------|--------------------|------------------|
| Agent | Reg↓ | NW↑ | NW ⁺ ↑ | UW↑ | EF1↑ |
| sonnet-3.7-c-5 | 29.49±.39 | 74.42±.10 | 52.46±.19 | 82.32±.06 | 31.40±.19 |
| sonnet-3.7-c-6 | 15.45±.47 | 74.51±.09 | 54.72±.20 | 82.88±.05 | 31.72±.21 |
| sonnet-3.7-r-c-0 | 24.19±.41 | 76.10±.08 | 52.85±.19 | 83.46±.07 | 41.53±.24 |
| gem-2.0-f-c-2 | 20.78±.33 | 76.59±.10 | 51.38±.20 | 83.46±.05 | 36.85±.21 |
| gem-2.0-f-c-5 | 22.28±.42 | 73.61±.14 | 50.18±.18 | 82.19±.06 | 30.94±.22 |
| 4o-c-4 | 15.63±.30 | 78.67±.07 | 53.50±.18 | 83.98±.05 | 36.04±.25 |
| 4o-c-5 | 18.97±.29 | 77.65±.10 | 53.44±.20 | 83.95±.06 | 36.47±.22 |
| 4o-c-6 | 15.75±.30 | 78.44±.10 | 48.53±.16 | 84.83±.07 | 33.79±.21 |
| o3-mini-c-0 | .84±.12** | 82.74±.08 | 60.97±.28 | 86.93±.06** | 54.22±.25 |
| nfsp | 18.35±.46 | 78.58±.10 | 55.71±.22 | 83.76±.07 | 55.73±.36 |
| R-NaD | 3.21±.25** | 80.56±.09 | 55.82±.19 | 85.94±.09** | 52.64±.31 |
| aspire | 21.00±.37 | 81.24±.09 | 58.54±.19 | 85.74±.06 | 44.10±.20 |
| soft | 29.28±.39 | 69.35±.10 | 44.92±.18 | 78.67±.04 | 27.31±.17 |
| tough | 49.52±.66 | 61.73±.13 | 6.49±.05 | 72.64±.04 | .70±.02 |
| walk | 51.85±.36 | 67.97±.01 | 0.00±.00 | 75.03±.00 | N/A |

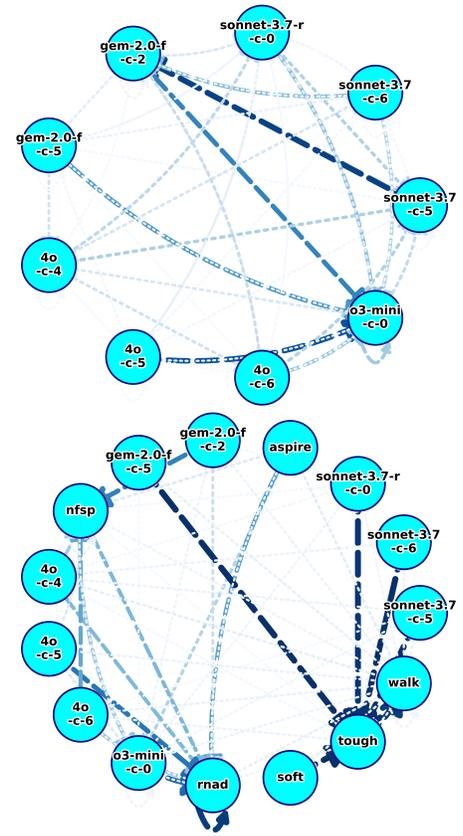


Figure 1: Left: Bootstrap means (\pm standard errors) for regret, normalized Nash and utilitarian welfare, and EF1 frequency (last three as percentages) in our LLM-only and LLM+baselines analyses on BGS0; best values are bold. Asterisks on o3-mini-c-0 and the next-ranked agent mark paired bootstrap significance: * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. Right: corresponding empirical best-response graphs (top: LLM-only; bottom: LLM+baselines). EF1 frequency is undefined for walk since its negotiations never end in ACCEPT.

meta-game (details in Appendix K.1). Figure 2 presents the provider-level BR graph.



Figure 2: Provider-level empirical best-response graph over BG1-BG6: Gemini is most often the best response to Anthropic, OpenAI to Gemini, and OpenAI to itself.

The patterns we find align with equilibrium behavior. Anthropic offers the most utility to opponents (in its own valuation), followed by OpenAI, then Gemini (Appendix K). The overall evidence supports the following account: Anthropic is generous yet more prone to accepting dominated offers, particularly against tough. Gemini adopts a tougher posture and exploits this, so it is frequently the best response to Anthropic. OpenAI is more robust to hardline bargaining while remaining reasonably generous, which makes it a strong best response to Gemini and often to itself. In aggregate, this produces provider-specific patterns with a non-transitive cycle of strengths and weaknesses.

7 Conclusion

We performed an in-depth assessment of LLM bargaining as a case study in evaluating the strategic capabilities of LLM agents. A key finding is that competition and cooperation were basically aligned in this mixed-motive setting. OpenAI models, especially o3-mini, performed best on both individual and social criteria, and were robust to inclusion of non-LLM baseline strategies. Gemini and especially Claude Sonnet were vulnerable to the hardline tough strategy, suggesting that excessive accommodation can be exploited by rigid opponents. Provider-level aggregation surfaced additional, non-transitive patterns among LLM families.

Given the speed of LLM advancement, findings about specific models have a short shelf-life. The more important results from this study are the methodological concepts developed and exercised in the process. By revealing qualitative strategic patterns and unsuspected vulnerabilities, we demonstrated that empirical meta-game analysis provides a practical lens for evaluating multiagent LLM behavior.

References

- Sahar Abdelnabi, Amr Gomaa, Sarath Sivaprasad, Lea Schönherr, and Mario Fritz. Cooperation, competition, and maliciousness: LLM-stakeholders interactive negotiation. In *38th Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2024. Datasets and Benchmarks Track.
- Gati V. Aher, Rosa I. Arriaga, and Adam Tauman Kalai. Using large language models to simulate multiple humans and replicate human subject studies. In *40th International Conference on Machine Learning (ICML)*, pages 337–371, 2023.
- Rohan Anil et al. Gemini: A family of highly capable multi-modal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Anthropic. The Claude 3 model family: Opus, Sonnet, Haiku. <https://www.anthropic.com/news/claude-3-family>, 2024. Model card, accessed 23 May 2025.
- Lisa P. Argyle, Ethan C. Busby, Nancy Fulda, Joshua R. Gubler, Christopher Rytting, and David Wingate. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351, 2023.
- David Balduzzi, Karl Tuyls, Julien Pérolat, and Thore Graepel. Re-evaluating evaluation. In *32nd Annual Conference on Neural Information Processing Systems (NeurIPS)*, pages 3268–3279, 2018.
- Federico Bianchi, Patrick John Chia, Mert Yuksekgonul, Jacopo Tagliabue, Dan Jurafsky, and James Zou. How well can LLMs negotiate? NegotiationArena platform and analysis. In *41st International Conference on Machine Learning (ICML)*, 2024.
- Eric Budish. The combinatorial assignment problem: Approximate competitive equilibrium from equal incomes. *Journal of Political Economy*, 119(6):1061–1103, 2011.
- Kris Cao, Angeliki Lazaridou, Marc Lanctot, Joel Z. Leibo, Karl Tuyls, and Stephen Clark. Emergent communication through negotiation. In *6th International Conference on Learning Representations (ICLR)*, 2018.
- Ioannis Caragiannis, David Kurokawa, Hervé Moulin, Ariel D. Procaccia, Nisarg Shah, and Junxing Wang. The unreasonable fairness of maximum Nash welfare. *ACM Transactions on Economics and Computation*, 7(3):12:1–12:32, 2019.
- Kushal Chawla, Ian Wu, Yu Rong, Gale Lucas, and Jonathan Gratch. Be selfish, but wisely: Investigating the impact of agent personality in mixed-motive human-agent interactions. In *Conference on Empirical Methods in Natural Language Processing*, pages 13078–13092, 2023.
- Jiangjie Chen, Siyu Yuan, Rong Ye, Bodhisattwa Prasad Majumder, and Kyle Richardson. Put your money where your mouth is: Evaluating strategic planning and execution of LLM agents in an auction arena. *arXiv preprint arXiv:2310.05746*, 2023.
- Yann Chevaleyre, Ulle Endriss, and Nicolas Maudet. Distributed fair allocation of indivisible goods. *Artificial Intelligence*, 242:1–22, 2017.
- Anthony Christopher Davison and David Victor Hinkley. *Bootstrap Methods and Their Application*. Cambridge University Press, 1997.
- Jinhao Duan, Renming Zhang, James Diffenderfer, Bhavya Kailkhura, Lichao Sun, Elias Stengel-Eskin, Mohit Bansal, Tianlong Chen, and Kaidi Xu. Gtbench: Uncovering the strategic reasoning limitations of LLMs via game-theoretic evaluations. In *38th Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2024.
- Ulrich Endriss, Nicolas Maudet, Fariba Sadri, and Francesca Toni. Negotiating socially optimal allocations of resources. *Journal of Artificial Intelligence Research*, 25:315–348, 2006.
- Caoyun Fan, Jindou Chen, Yaohui Jin, and Hao He. Can large language models serve as rational players in game theory? A systematic analysis. In *38th AAAI Conference on Artificial Intelligence (AAAI)*, pages 17960–17967, 2024.
- Roger Fisher, William L. Ury, and Bruce Patton. *Getting to Yes: Negotiating Agreement without Giving In*. Penguin Books, 3rd edition, 2011.
- John C. Harsanyi. On the rationality postulates underlying the theory of cooperative games. *Journal of Conflict Resolution*, 5(2):179–196, 1961.
- He He, Derek Chen, Anusha Balakrishnan, and Percy Liang. Decoupling strategy and generation in negotiation dialogues. In *Conference on Empirical Methods in Natural Language Processing*, pages 2333–2343, 2018.
- Johannes Heinrich and David Silver. Deep reinforcement learning from self-play in imperfect-information games. *arXiv preprint arXiv:1603.01121*, 2016.
- Johannes Heinrich, Marc Lanctot, and David Silver. Fictitious self-play in extensive-form games. In *32nd International Conference on Machine Learning (ICML)*, pages 805–813, 2015.
- Daniel Hennes, Dustin Morrill, Shayegan Omidshafiei, et al. Neural replicator dynamics: Multiagent learning via hedging policy gradients. In *19th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, pages 492–501, 2020.
- John J. Horton. Large language models as simulated economic agents: What can we learn from homo silicus? Technical report, National Bureau of Economic Research, 2023.
- Hengyuan Hu, Adam Lerer, Alex Peysakhovich, and Jakob Foerster. “Other-play” for zero-shot coordination. In *37th International Conference on Machine Learning (ICML)*, pages 4399–4410, 2020.
- Wenyue Hua et al. Game-theoretic LLM: Agent workflow for negotiation games. *arXiv preprint arXiv:2411.05990*, 2024.
- Wenyue Hua, Ollie Liu, Lingyao Li, Alfonso Amayuelas, Julie Chen, Lucas Jiang, Mingyu Jin, Lizhou Fan, Fei Sun, William Wang, Xintong Wang, and Yongfeng Zhang. Game-theoretic LLM: Agent workflow for negotiation games. *arXiv preprint arXiv:2411.05990*, 2024.

- Aaron Hurst et al. GPT-4o system card. <https://cdn.openai.com/papers/gpt-4o-system-card.pdf>, 2024. Last accessed: May 23, 2025.
- Patrick R. Jordan, Christopher Kiekintveld, and Michael P. Wellman. Empirical game-theoretic analysis of the TAC supply chain game. In *6th International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, pages 1188–1195, 2007.
- Scott M. Jordan, Yash Chandak, Daniel Cohen, Mengxue Zhang, and Philip S. Thomas. Evaluating the performance of reinforcement learning algorithms. In *37th International Conference on Machine Learning (ICML)*, pages 4962–4973, 2020.
- Christine Konicki, Mithun Chakraborty, and Michael P. Wellman. Policy abstraction and Nash refinement in tree-exploiting PSRO. In *24th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, pages 1163–1171, 2025.
- Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, et al. Openspiel: A framework for reinforcement learning in games. *arXiv preprint arXiv:1908.09453*, 2019.
- Mike Lewis, Denis Yarats, Yann N. Dauphin, Devi Parikh, and Dhruv Batra. Deal or no deal? End-to-end learning for negotiation dialogues. In *Conference on Empirical Methods in Natural Language Processing*, 2017.
- Zun Li and Michael P. Wellman. A meta-game evaluation framework for deep multiagent reinforcement learning. In *33rd International Joint Conference on Artificial Intelligence (IJCAI)*, pages 148–156, 2024.
- Zun Li, Marc Lanctot, Kevin McKee, Luke Marris, Ian Gemp, Daniel Hennes, Kate Larson, Yoram Bachrach, and Michael P. Wellman. Combining deep reinforcement learning and search with generative models for game-theoretic opponent modeling. In *34th International Joint Conference on Artificial Intelligence (IJCAI)*, 2025.
- Richard J. Lipton, Evangelos Markakis, Elchanan Mossel, and Amin Saberi. On approximately fair allocations of indivisible goods. In *5th ACM Conference on Electronic Commerce (EC)*, pages 125–131, 2004.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. AgentBench: Evaluating LLMs as agents. In *12th International Conference on Learning Representations (ICLR)*, 2024.
- Siqi Liu, Ian Gemp, Luke Marris, Georgios Piliouras, Nicolas Heess, and Marc Lanctot. Re-evaluating open-ended evaluation of large language models. In *13th International Conference on Learning Representations (ICLR)*, 2025.
- Panayotis Mertikopoulos, Christos H. Papadimitriou, and Georgios Piliouras. Cycles in adversarial regularized learning. In *29th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 2703–2717, 2018.
- Volodymyr Mnih et al. Human-level control through deep reinforcement learning. *Nature*, 518:529–533, 2015.
- Yasser Mohammad, Shinji Nakadai, and Amy Greenwald. NegMAS: A platform for automated negotiations. In *International Conference on Principles and Practice of Multi-Agent Systems (PRIMA)*, volume 12408 of *Lecture Notes in Computer Science*, pages 403–419. Springer, 2020.
- OpenAI. OpenAI o3 and o4-mini system card. <https://cdn.openai.com/papers/o3-o4-mini-system-card.pdf>, April 2025. Accessed 23 May 2025.
- Georgios Papoudakis, Filippos Christianos, Lukas Schäfer, and Stefano V. Albrecht. Benchmarking multi-agent deep reinforcement learning algorithms in cooperative tasks. In *35th Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2021. Datasets and Benchmarks Track.
- Julien Perolat et al. From Poincaré recurrence to convergence in imperfect information games: Finding equilibrium via regularization. pages 8525–8535, 2021.
- Julien Perolat et al. Mastering the game of Stratego with model-free multiagent reinforcement learning. *Science*, 378(6623):990–996, 2022.
- Giorgio Piatti, Zhijing Jin, Max Kleiman-Weiner, Bernhard Schölkopf, Mrinmaya Sachan, and Rada Mihalcea. Cooperate or collapse: Emergence of sustainable cooperation in a society of LLM agents. In *38th Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2024.
- Sara Ramezani and Ulle Endriss. Nash social welfare in multiagent resource allocation. In *Agent-Mediated Electronic Commerce: Designing Trading Strategies and Mechanisms for Electronic Markets*, volume 59 of *Lecture Notes in Business Information Processing*, pages 117–131. Springer, 2010.
- Arunesh Sinha and Michael P. Wellman. Incentivizing collaboration in a competition. In *18th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, pages 556–564, 2019.
- Constantino Tsallis. Possible generalization of Boltzmann-Gibbs statistics. *Journal of Statistical Physics*, 52(1-2):479–487, 1988.
- Jeffrey S. Vitter. Random sampling with a reservoir. *ACM Transactions on Mathematical Software*, 11(1):37–57, March 1985.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *36th Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- Michael P. Wellman, Karl Tuyls, and Amy Greenwald. Empirical game-theoretic analysis: A survey. *Journal of Artificial Intelligence Research*, 82:1017–1076, 2025.
- Bryce Wiedenbeck, Ben-Alexander Cassell, and Michael P. Wellman. Bootstrap statistics for empirical games. In *13th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, pages 597–604, 2014.

Submission 1298: APPENDICES

A Experimental setup details

| Game class | T | q | Game | γ | R | #settings in cohort | Total #settings |
|------------|-----|---|------|----------|-----|---------------------|-----------------|
| BGL | 5 | $\forall t \in \{1, \dots, 5\}. q_t \sim \text{Poisson}(4)$ | BG1 | 0.9 | 3 | 11,630 | 27,413 |
| | | | BG2 | 0.98 | 3 | 7,933 | |
| | | | BG3 | 0.98 | 5 | 7,850 | |
| BGS | 3 | (7, 4, 1) | BG4 | 0.9 | 3 | 25,700 | 75,740 |
| | | | BG5 | 0.98 | 3 | 24,900 | |
| | | | BG6 | 0.98 | 5 | 25,100 | |

Table 1: Statistics of game instances generated for our bootstrapped empirical meta-game analysis, as described in Section 5. T denotes the number of item types in the pool; q is the T -length vector of the numbers of items of all types in the pool; γ is the (exponential) discount rate shared by both players; R is the number of bargaining rounds in an instance.

B Negotiation mistakes

In any bargaining interaction, certain moves should be avoided to maintain rational play. The five key mistakes, with illustrative examples, are:

- M1: Making an offer worse than your previous offer.** You reject an opponent’s proposal that would have given you more value than the counteroffer you then make.
Example: Suppose your private values per unit are $v = [4, 2]$ for items A and B. Your opponent offers you (3, 0), worth $3 \cdot 4 + 0 \cdot 2 = 12$. You then counter with (2, 0) (worth 8), you have made a strictly worse proposal than the one you just rejected.
- M2: Making an offer worse for you than your outside offer.** You propose a split in which the items you keep are worth less than your guaranteed fallback (outside offer).
Example: Your outside-offer value is 10. You propose to keep (1, 1), worth $1 \cdot 4 + 1 \cdot 2 = 6$, even though walking away would guarantee you 10.
- M3: Offering no items or all items.** Issuing an extreme division—giving everything away or keeping everything.
Example: In Round 2 of a 5-round game, you offer to give the opponent all items (0, 0) for yourself.
- M4: Accepting an offer worse for you than your outside offer.** You agree to a deal that yields you less value than you would have secured by walking away to your outside option.
Example: Your outside-offer value is 8. Your opponent offers you (1, 0), worth 4. If you accept, you receive $4 < 8$, which is dominated by simply walking away.
- M5: Walking away from an offer better than your outside offer.** You reject a valid split that provides you more value than your fallback, thus foregoing a strictly superior outcome.
Example: Your outside-offer value is 9. The opponent offers you (2, 1), worth $2 \cdot 4 + 1 \cdot 2 = 10$. Walking away would leave you with 9, so walking away is strictly dominated by the current offer on the table.

C Prompt engineering for negotiation

This section describes each incremental "circle" of prompts provided to the negotiation models. Each circle adds guidance or strategic considerations aimed at progressively improving the reasoning and negotiation outcomes. The circles follow a progression, with each level introducing more sophisticated negotiation concepts.

The negotiation history is maintained throughout the game, displaying each player’s actions chronologically. The current offer on the table is explicitly stated, and the action prompt varies based on whether the agent is Player 1 making the first move, has a current offer to consider, or needs to make a counteroffer.

Prompts for the 3 item games are modified in accordance with the differing item count in its games.

C.1 System Prompt

All LLM agents receive a consistent system prompt that establishes their role as negotiators and enforces the required response format:

System Prompt

You are an AI negotiator participating in a negotiation game.
You must respond in one of these formats *only*:

1. {"action": "ACCEPT"}
2. {"action": "WALK"}
3. {"action": "COUNTEROFFER", "offer": [n1, n2, ...]}

Ensure your response is valid JSON and matches one of these exact formats.

This system prompt is sent with every API call across all models and prompt circles, ensuring consistent response formatting. The detailed game instructions and strategic guidance are provided in the user message, which varies according to the prompt circle level described below.

C.2 Prompting "Circles"

Circle 0: The Vestibule (Setting the Stage). The baseline negotiation scenario provides the foundational rules of the negotiation game, including player identities, item quantities, private valuations, outside offers, and negotiation structure. No additional guidance is provided beyond the basic game mechanics.

Circle 0 Prompt

You and another agent have to negotiate a division of items between the two of you. You are Player {my_player_num} and the other agent is Player {other_player_num}. There are T types of items, called item 1 through item T. There are {quantities} to divide. Both you and Player {other_player_num} have a private value per unit of each item type. These values are drawn from a uniform random distribution, ranging from 1 to V-1. Your private values are {values}. You have a private outside offer drawn from a uniform random distribution ranging from {p1_outside_offer[0]} to your total value of all items, which is {p1_outside_offer[1]}. Player {other_player_num} has a private outside offer drawn from a uniform random distribution ranging from 1 to their total value of all items. Your outside offer value is {w}.

The negotiation proceeds in {R} rounds. There is a discount rate $\gamma = \{g\}$, such that if the process concludes after r rounds the overall value of the negotiation to each player is their value for the outcome multiplied by γ to the power (r-1). At each round, Player 1 takes an action, followed by Player 2. The possible actions are to ACCEPT the other player's current offer (if any), make a COUNTEROFFER, or WALK away. If the game gets to the last round, and player 2 chooses to make a counteroffer, this is treated as a WALK. If a player chooses ACCEPT, the negotiation ends in a deal to divide the items according to the accepted offer. The value of an outcome is determined by each player's private values per unit of each item and the quantities they receive in the deal. This value is adjusted by the discount factor, which is used to compute the present value of the negotiation outcome. If a player chooses WALK, the negotiation ends without a deal, and each player receives the value of their private outside offer.

Circle 1: The Limbo of Objectives. Circle 1 explicitly introduces the negotiation objective, reminding the agent to maximize its negotiation outcome and to carefully consider its guaranteed alternative, the outside offer.

Circle 1 Prompt

[All instructions from Circle 0, plus:]
Your objective is to maximize your value of the outcome of the negotiation game.
Remember, you have a guaranteed alternative: your outside offer.

Circle 2: The Circle of Calculations. Circle 2 provides concrete numerical examples illustrating how agents should evaluate potential counteroffers by explicitly calculating their total value, ensuring the counteroffers are superior to their outside offers.

Circle 2 Prompt

[All instructions from Circle 1, plus:]
Before making any counteroffer, you should calculate its total value to you and compare it to your outside offer value of w . For example, if you were considering offering the other player 2 units of each item (keeping 3 units of each for yourself), you would calculate: 3 units of item 1 = $3 \times \text{values}[0] = 3 \times \text{values}[0]$ (multiplying units by your value per unit) 3 units of item 2 = $3 \times \text{values}[1] = 3 \times \text{values}[1]$ (multiplying units by your value per unit) 3 units of item 3 = $3 \times \text{values}[2] = 3 \times \text{values}[2]$ (multiplying units by your value per unit) 3 units of item 4 = $3 \times \text{values}[3] = 3 \times \text{values}[3]$ (multiplying units by your value per unit) 3 units of item 5 = $3 \times \text{values}[4] = 3 \times \text{values}[4]$ (multiplying units by your value per unit)
Total value = $\text{sum}([3 \times \text{values}[i] \text{ for } i \text{ in range}(T)])$ (sum of all item values)

Circle 3: The Circle of Systematic Analysis. Circle 3 introduces a systematic three-step questioning strategy designed to facilitate a comprehensive analysis:

1. Analysis of the current negotiation situation
2. Assessment of the value of current and potential offers
3. Decision-making based on comparative analyses

Circle 3 Prompt

[All instructions from Circle 2, plus:]
The following step-by-step questions are designed to guide you through a comprehensive analysis. By systematically addressing these questions, you can evaluate the current state of the negotiation, assess potential offers, and make informed decisions. You must use the information that you acquired through the step-by-step questioning above to decide what action you will make.
Let's walk through this step by step:

1. First, analyze the current situation:
 - What is my outside offer value?
 - What are the values of the items involved?
 - What is the total pool of items?
 - How does the discount factor influence the value of accepting the current offer versus waiting for future offers?
2. Assess the value of offers:
 - For the current offer (if any): What is my total value if I accept it?
 - For potential counteroffers: What would be my total value for different proposed divisions?
 - How do these values compare to my outside offer value?
3. Make a decision based on the analysis:
 - Should I accept the current offer?
 - Should I walk away and take my outside offer?
 - Or should I propose a specific counteroffer?

Circle 4: The Circle of Common Errors. Circle 4 identifies five common negotiation mistakes that conflict with the negotiation objectives:

- **Mistake 1:** Making an offer worse than your previous offer. This occurs when you reject an offer better for you than the one you subsequently propose.
- **Mistake 2:** Making an offer worse for you than your outside offer. This happens if you propose giving away so much that what you keep is worth less than your guaranteed alternative.
- **Mistake 3:** Offering no items or all items. Offering nothing (or everything) to the opponent (in the early or middle rounds) can be a clear suboptimal move.
- **Mistake 4:** Accepting an offer worse for you than your outside offer. This occurs if you accept a division that yields a payoff lower than your guaranteed fallback.

- **Mistake 5:** Walking away from an offer better than your outside offer. This occurs when you reject a division that actually yields a higher payoff than your fallback.

Circle 4 Prompt

[All instructions from Circle 3, plus:]

In the bargaining game, there are five mistakes you can make that conflict with your objectives. While these aren't the only possible errors, they represent undesirable negotiation behaviors that can undermine your payoff or cause you to miss out on better deals.

1. **Mistake 1:** Making an offer worse than your previous offer. This occurs when you reject an offer better for you than the one you subsequently propose.
2. **Mistake 2:** Making an offer worse for you than your outside offer. This happens if you propose giving away so much that what you keep is worth less than your guaranteed alternative, which is your outside offer.
3. **Mistake 3:** Offering no items or all items. Offering nothing (or everything) to the opponent (in the early or middle rounds) can be a clear sub-optimal move.
4. **Mistake 4:** Accepting an offer worse for you than your outside offer. This occurs if you accept a division that yields a payoff lower than your guaranteed fallback.
5. **Mistake 5:** Walking away from an offer better than your outside offer. This occurs when you reject a division that actually yields a higher payoff than your fallback.

Circle 5: The Circle of Error Prevention. Circle 5 provides explicit examples and strategic recommendations to prevent the identified mistakes, emphasizing the calculation of offer values and iterating proposals until surpassing the outside offer value. The specific example is dynamically generated for each player in each game instance. This circle introduces a concrete example of a suboptimal offer to illustrate the importance of value calculation.

Circle 5 Prompt

[All instructions from Circle 4, plus:]

To prevent these mistakes, adopt a strategy similar to the following example: before making any counter-offer, calculate its total value to you and compare it to your outside-offer value. For instance, suppose you keep only {example_offer_less_than_outside_offer_self} items and offer the rest to the other party. Your value would be:

$$\begin{aligned} & \{values[0]\} \times \{example_offer_less_than_outside_offer_self[0]\} + \\ & \{values[1]\} \times \{example_offer_less_than_outside_offer_self[1]\} + \\ & \{values[2]\} \times \{example_offer_less_than_outside_offer_self[2]\} + \\ & \{values[3]\} \times \{example_offer_less_than_outside_offer_self[3]\} + \\ & \{values[4]\} \times \{example_offer_less_than_outside_offer_self[4]\} \end{aligned}$$

which is $\{np.dot(values, example_offer_less_than_outside_offer_self)\}$ (sum of all item values) and is less than your outside offer of {w}. If your proposed offer results in a value lower than your outside offer, continue iterating until you develop a more advantageous offer that is better than your outside offer.

This reasoning can be applied to each of the five highlighted mistakes to ensure that your offers align with your objectives and avoid undesirable negotiation behaviors.

Circle 6: The Circle of Strategic Inference. Circle 6 introduces strategic inference about the opposing agent's valuations based on their negotiation behavior. Agents are guided to infer relative valuations from the distribution of items in the opponent's offers, leveraging this information to inform their own negotiation strategy.

Circle 6 Prompt

[All instructions from Circle 5, plus:]

Keep in mind the offers the opposing agent makes reflects its own values. If their offer includes most or all units of a particular item, it might indicate that the agent does not highly value that item, whereas offering none could suggest the opposite. You can use this kind of evidence to help inform your decision-making.

C.3 Response format and history tracking

All prompts include standardized action specification requirements and maintain negotiation history across rounds:

Response Format

Please show your reasoning step by step, then provide your action in one of these formats in your response (if you do not do this your response will be invalid):
{`"action": "ACCEPT"`} -- to accept the current offer
{`"action": "WALK"`} -- to walk away from negotiations
{`"action": "COUNTEROFFER", "offer": [n1, n2, ...]`} -- where n1, n2, ... are numbers representing the number of units of each item being offered to the other player as part of the counter-offer.

Any response not in these exact formats will be invalid and treated as a WALK. If you provide a counter-offer, it must be a valid offer; otherwise it will be treated as a WALK.

It is now round {r}.

Negotiation history: [History of all previous rounds, including offers made and actions taken]

Current offer on the table (the amount of each item being offered to you):

{`current_offer`}

What is your action? [Action prompt specific to the current game state]

D Circle selection methodology

We paired each of the LLMs GPT-4o, o3-mini, Anthropic Claude 3.7 Sonnet, and Gemini 2.0 Flash with each of our 7 prompt circles (as described in Section C) and sampled 100 self-play instances of BGL. For each instance, we recorded the number of times each mistake (as defined in Section B) was committed by both players over all rounds. For each LLM, we use a radar chart to visualize the cumulative frequencies of all 5 mistakes over all 100 instances for each circle in Fig. 3.

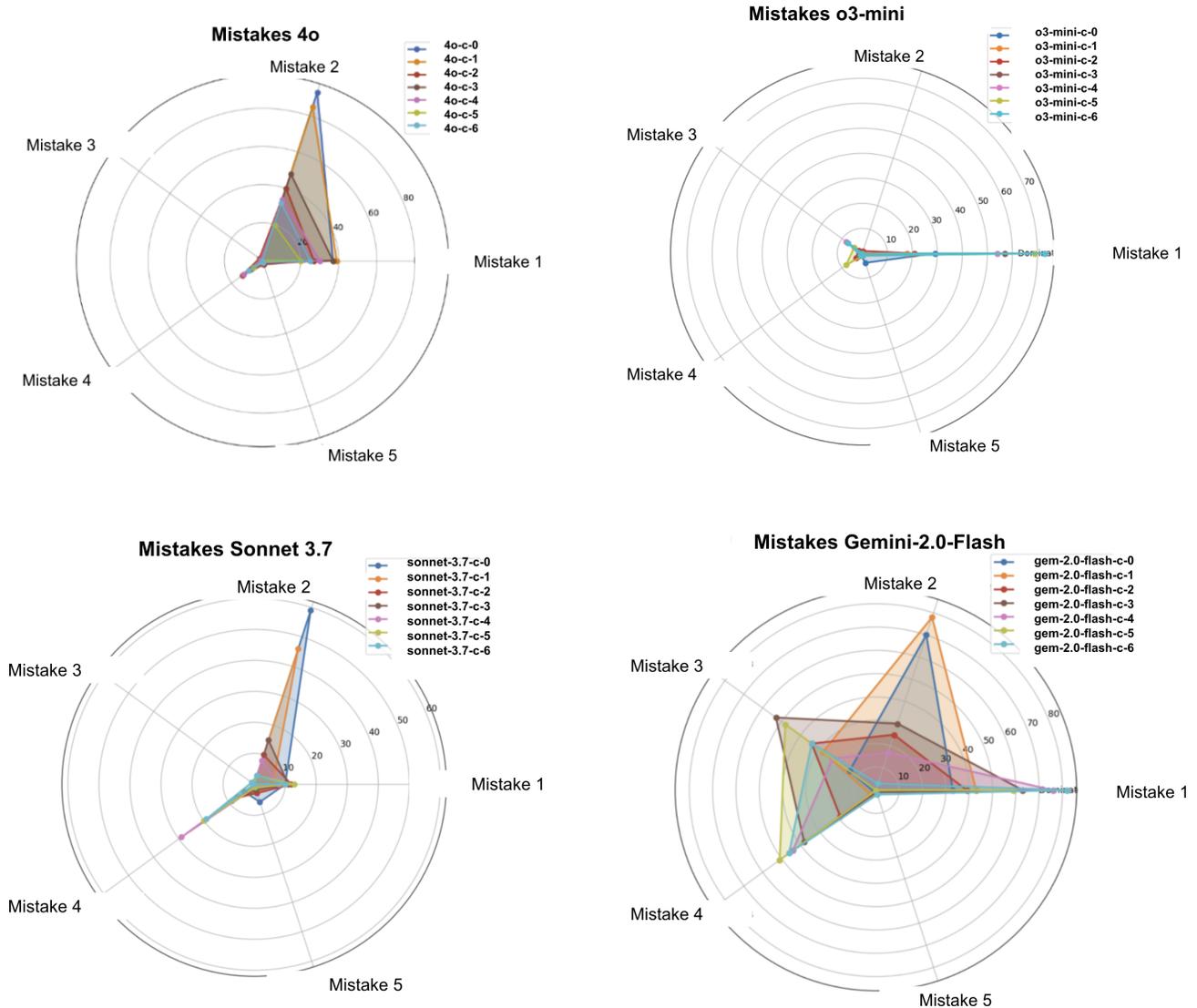


Figure 3: **Mistake counts for each negotiation strategy as a function of prompt "circle."** Each radar chart shows the number of times mistakes of each kind (as defined in Section B) were made by the corresponding LLM when paired with each incremental prompt circle (0 – 6) over all self-play histories. Lower counts indicate better adherence to rational negotiation behavior.

Although we crafted our prompt circles to provide progressively more elaborate guidance to the LLM to become a better negotiator (in particular, to avoid the potential mistakes we identified), our intuition was not generally borne out by our experimental results, so we finally had to make a judgment call to pick our final set LLM-circle pairings for our meta-game analysis. GPT-4o and Claude Sonnet 3.7 achieved a substantial reduction in mistake frequencies with higher circles, although not a consistent Pareto improvement over all mistakes with an increase in the circle number. Gemini 2.0 Flash displayed minimal sensitivity to prompt complexity. Most surprisingly, the performance OpenAI’s reasoning model o3-mini along most dimensions significantly deteriorated with higher circles, suggesting that additional instruction may interfere with their inherent reasoning capabilities. In light of this last finding, we did not subject Sonnet 3.7 Reasoning to the same treatment as the above

4 models and used its pairing with Circle 0 only.

Our final 9 LLM agents are:

- 4o-c-4,
- 4o-c-5,
- 4o-c-6,
- o3-mini-c-0,
- sonnet-3.7-c-5,
- sonnet-3.7-c-6,
- sonnet-3.7-r-c-0,
- gem-2.0-f-c-2,
- gem-2.0-f-c-5

E Empirical Game Average Max-Entropy Nash Equilibria

In this section, we report the average maximum-entropy Nash equilibrium computed on every cohort of bootstrap samples in LLM-only and LLM+baselines.

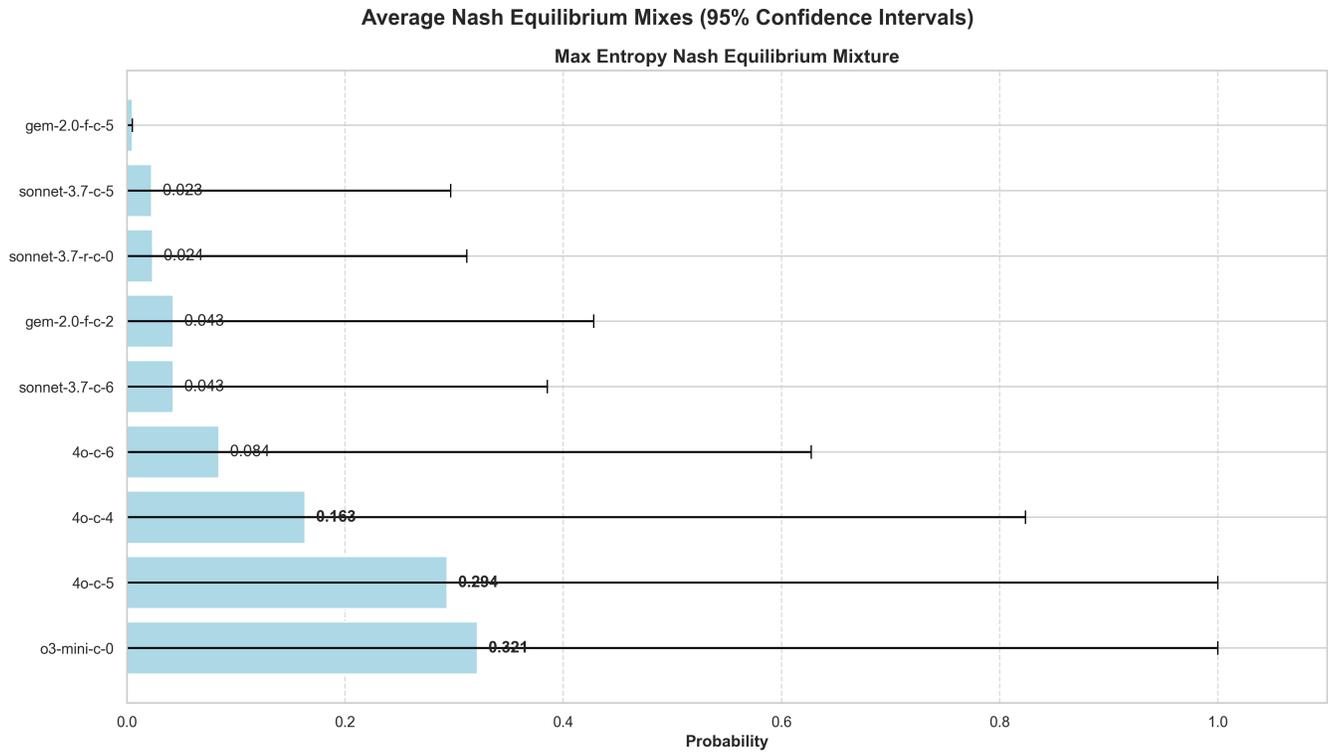


Figure 4: BG1, LLM-only

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

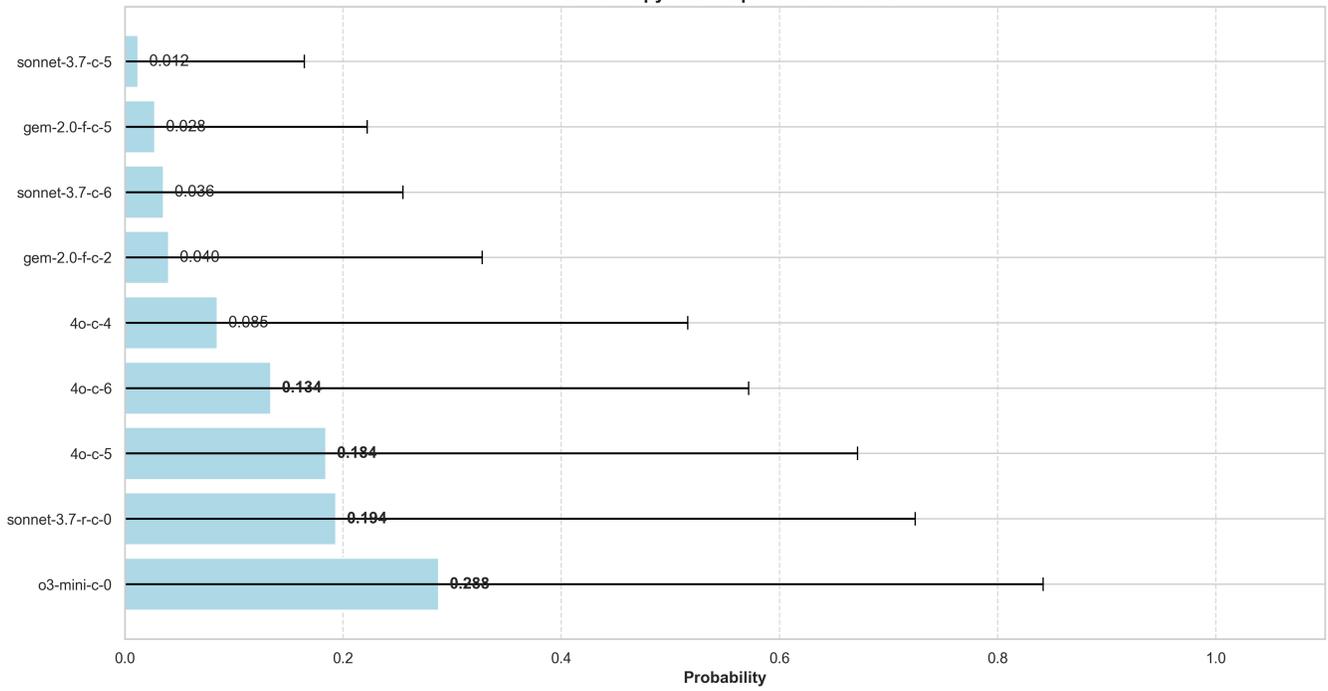


Figure 5: BG2, LLM-only

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

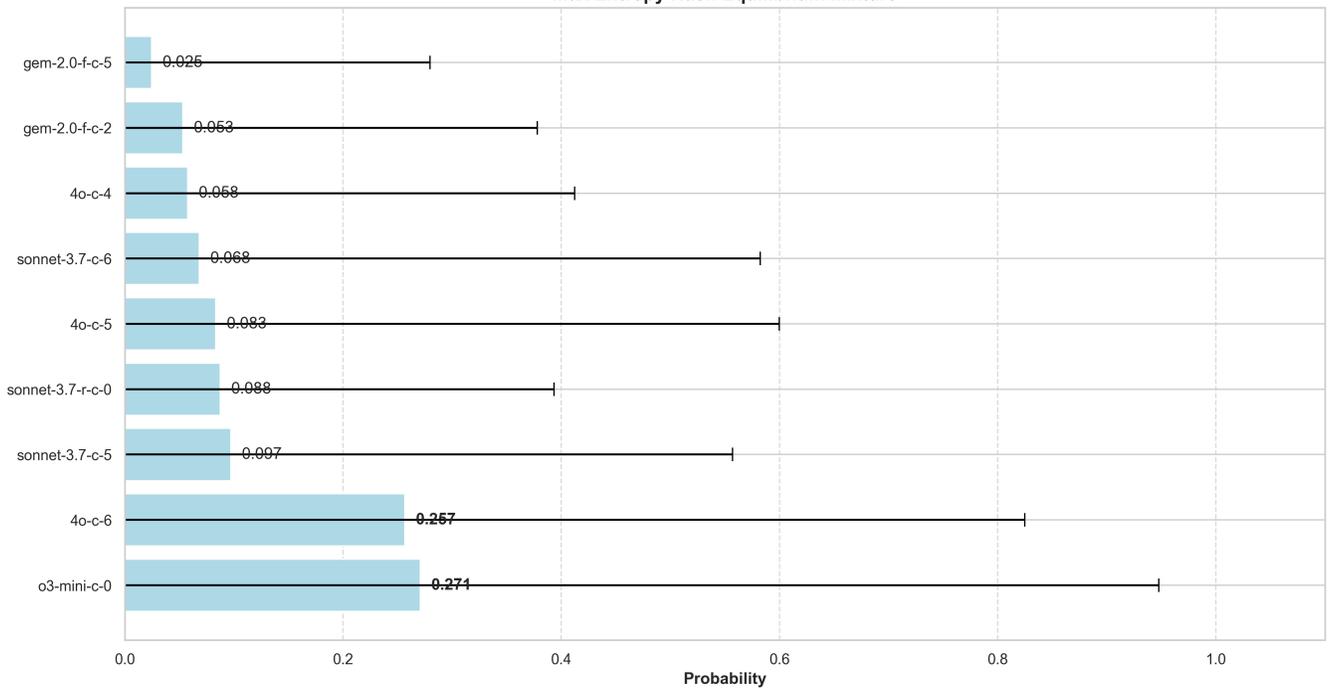


Figure 6: BG3, LLM-only

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

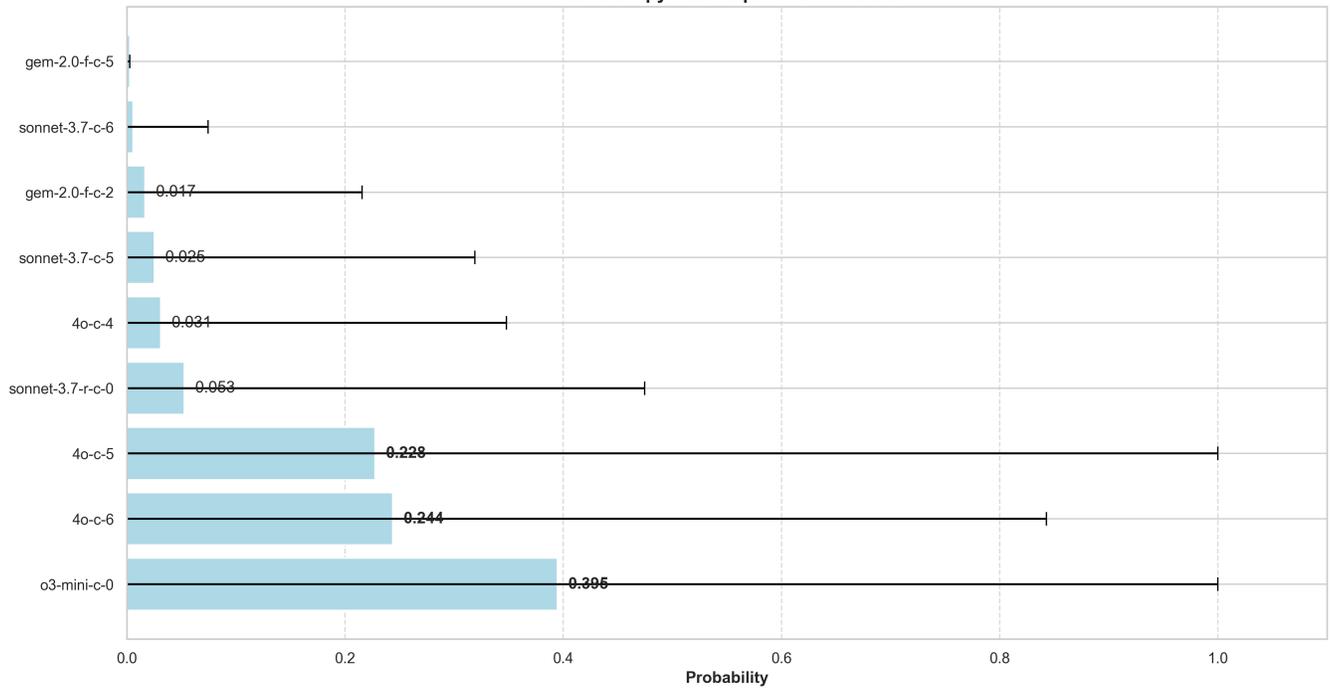


Figure 7: BGL0, LLM-only

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

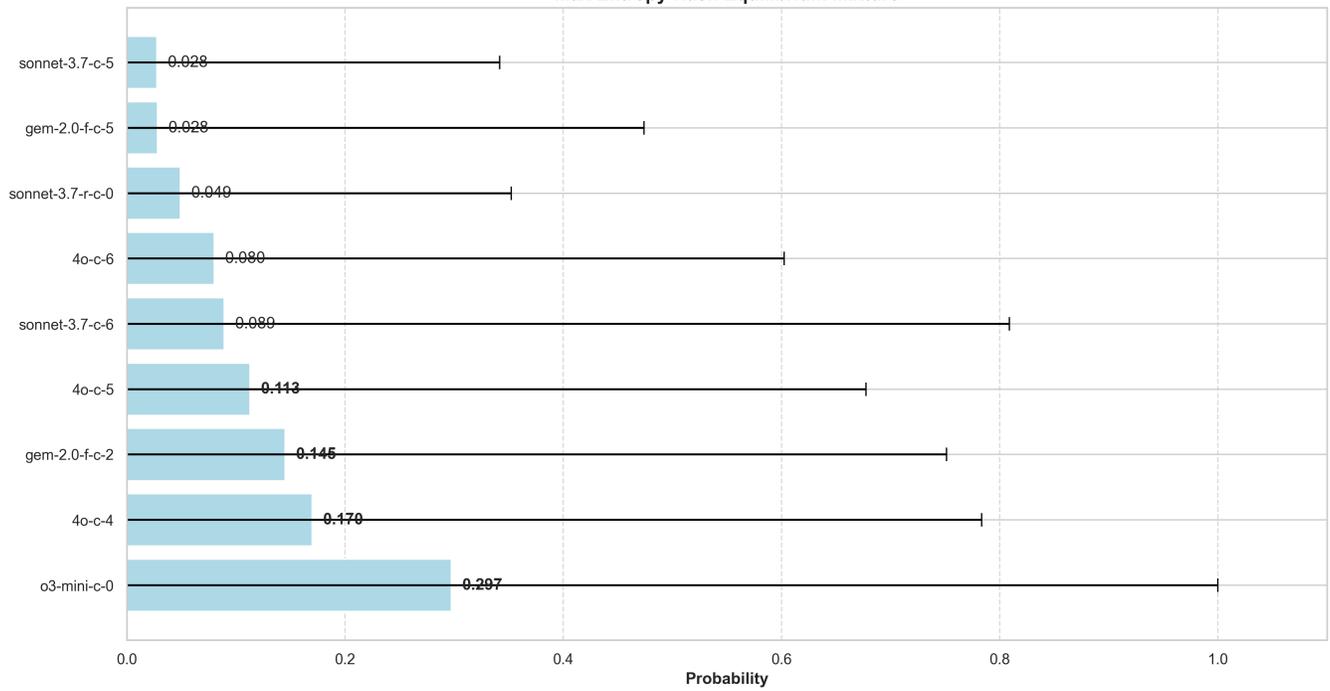


Figure 8: BG4, LLM-only

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

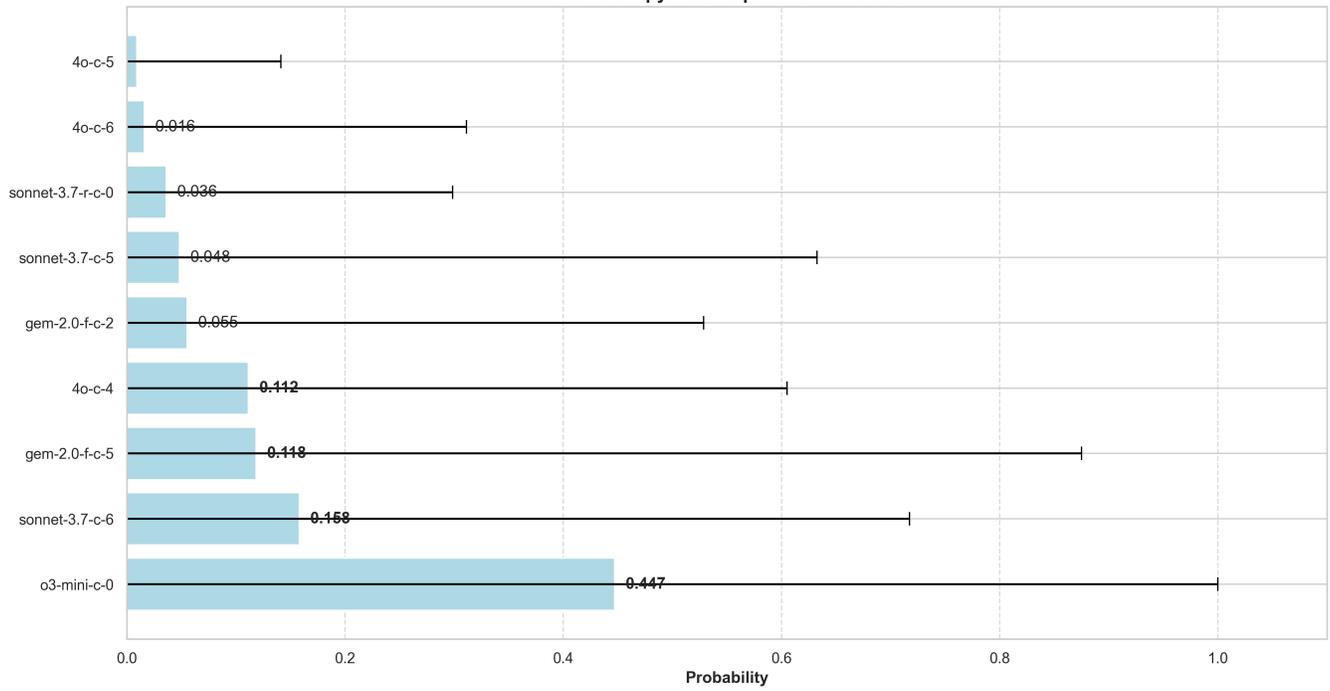


Figure 9: BG5, LLM-only

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

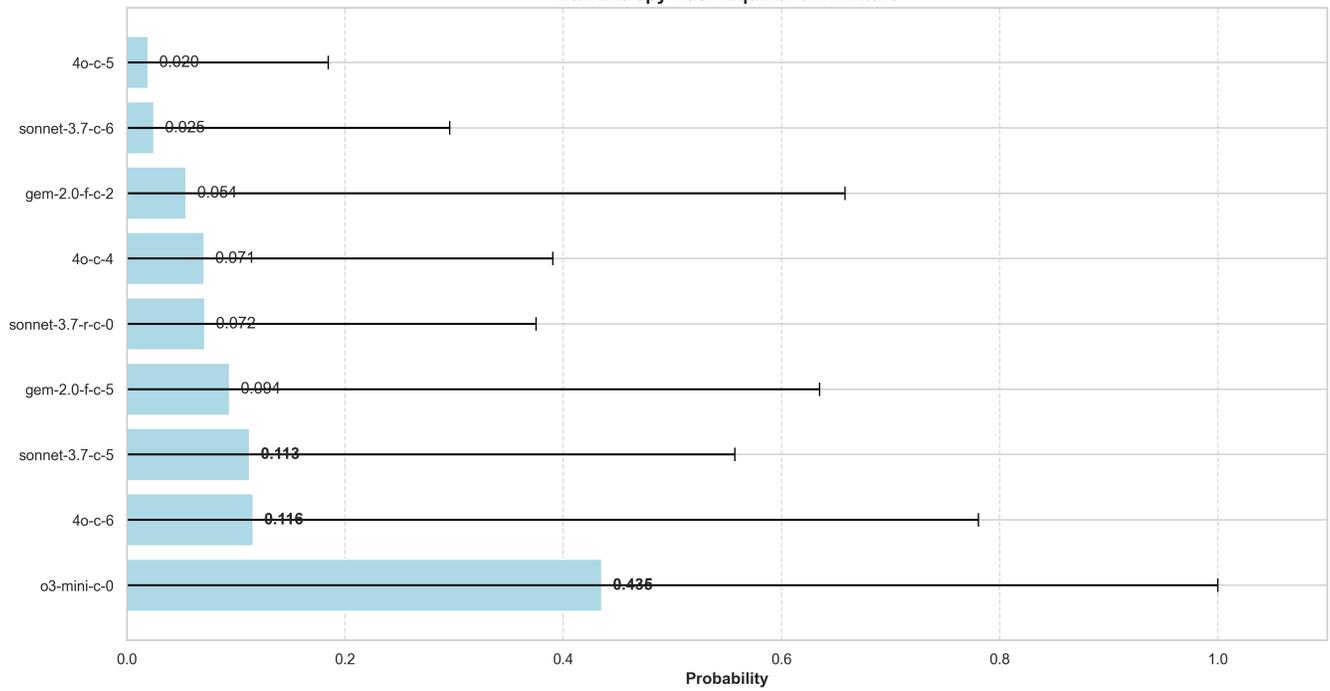


Figure 10: BG6, LLM-only

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

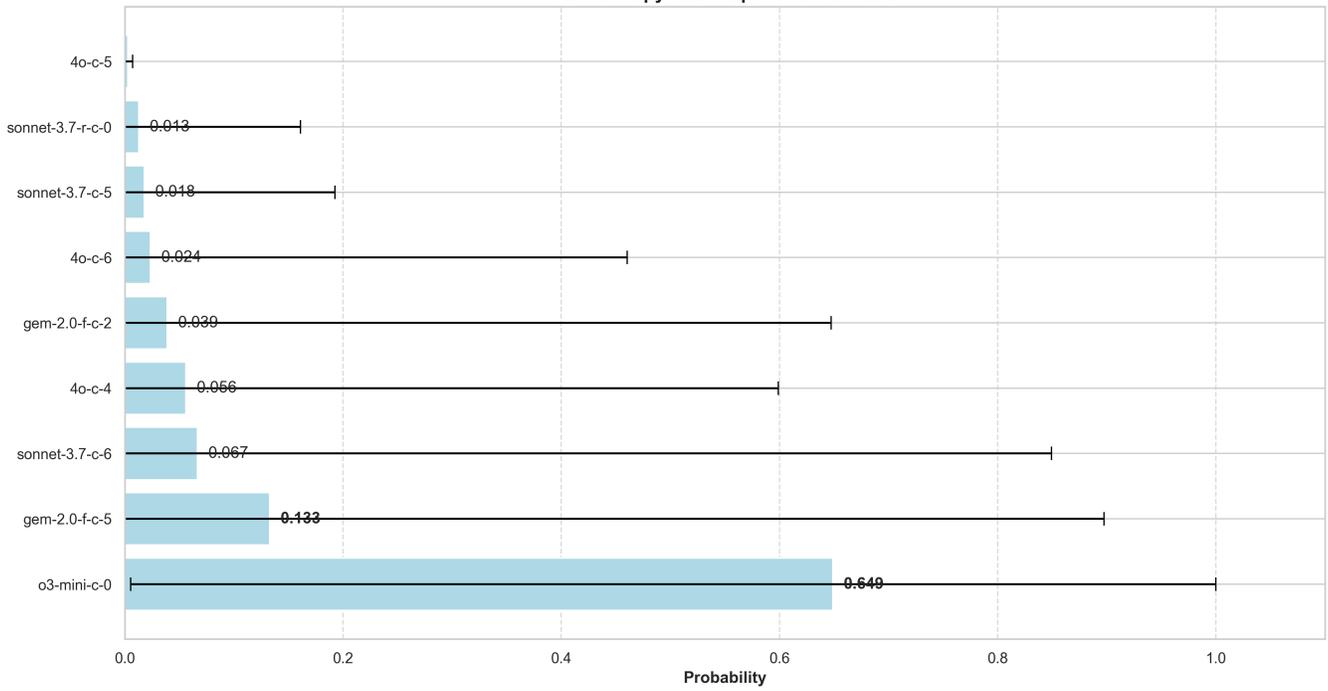


Figure 11: BGS0, LLM-only

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

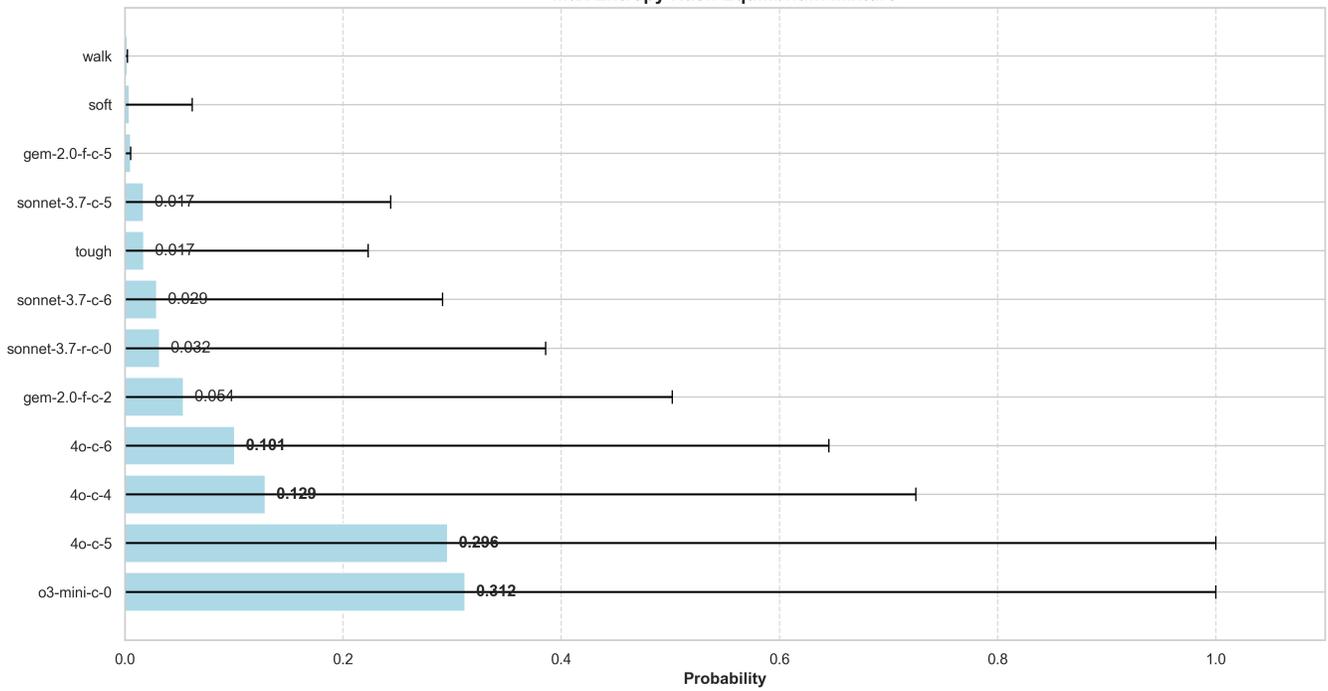


Figure 12: BG1, LLM+baselines

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

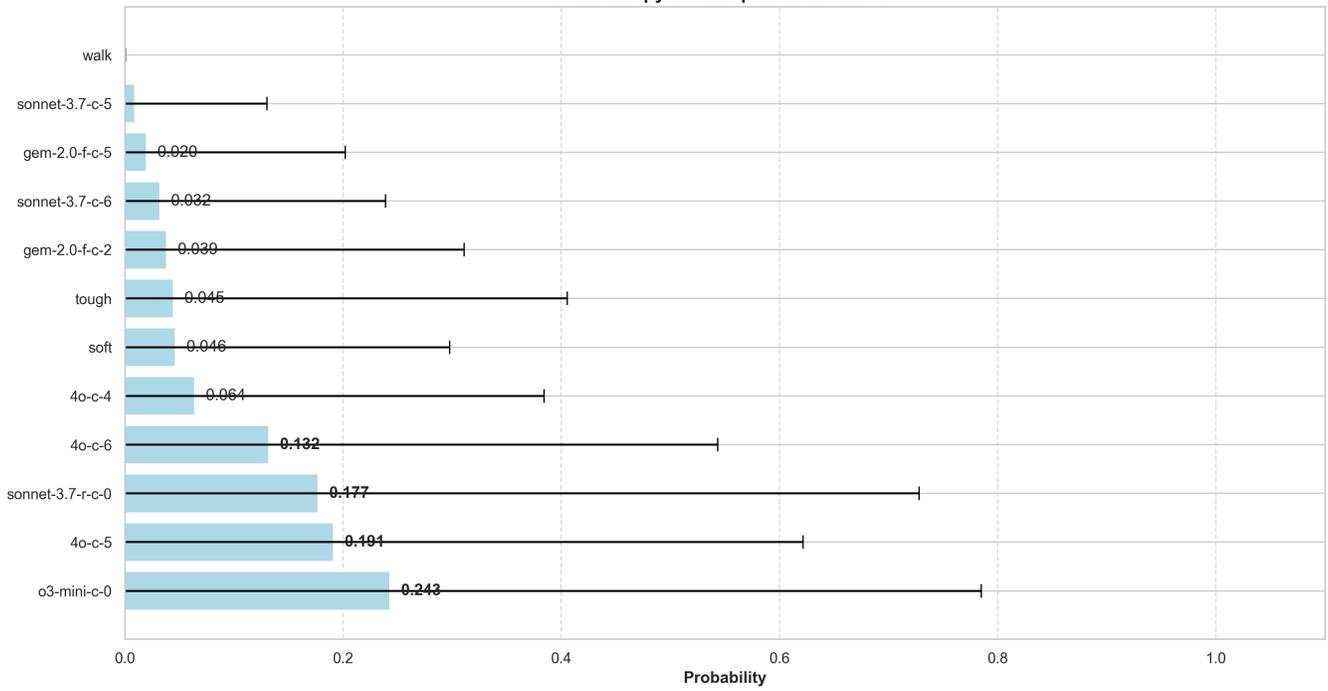


Figure 13: BG2, LLM+baselines

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

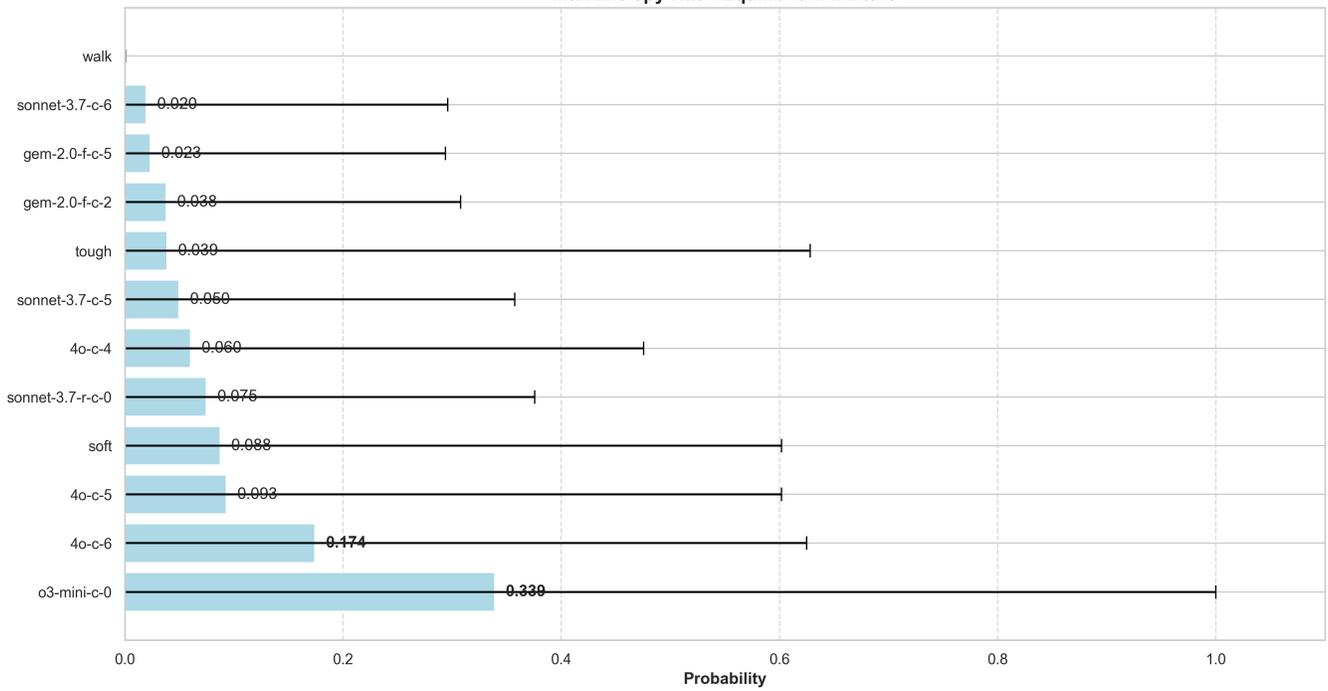


Figure 14: BG3, LLM+baselines

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

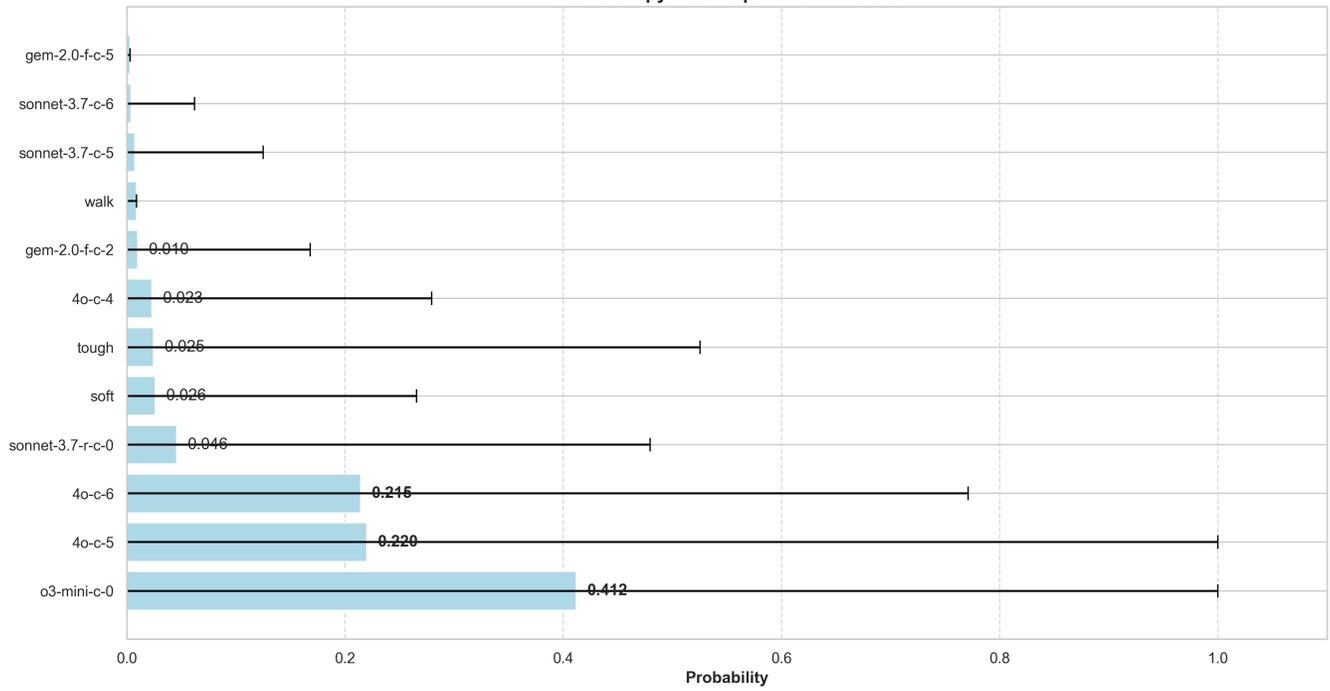


Figure 15: BGL0, LLM+baselines

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

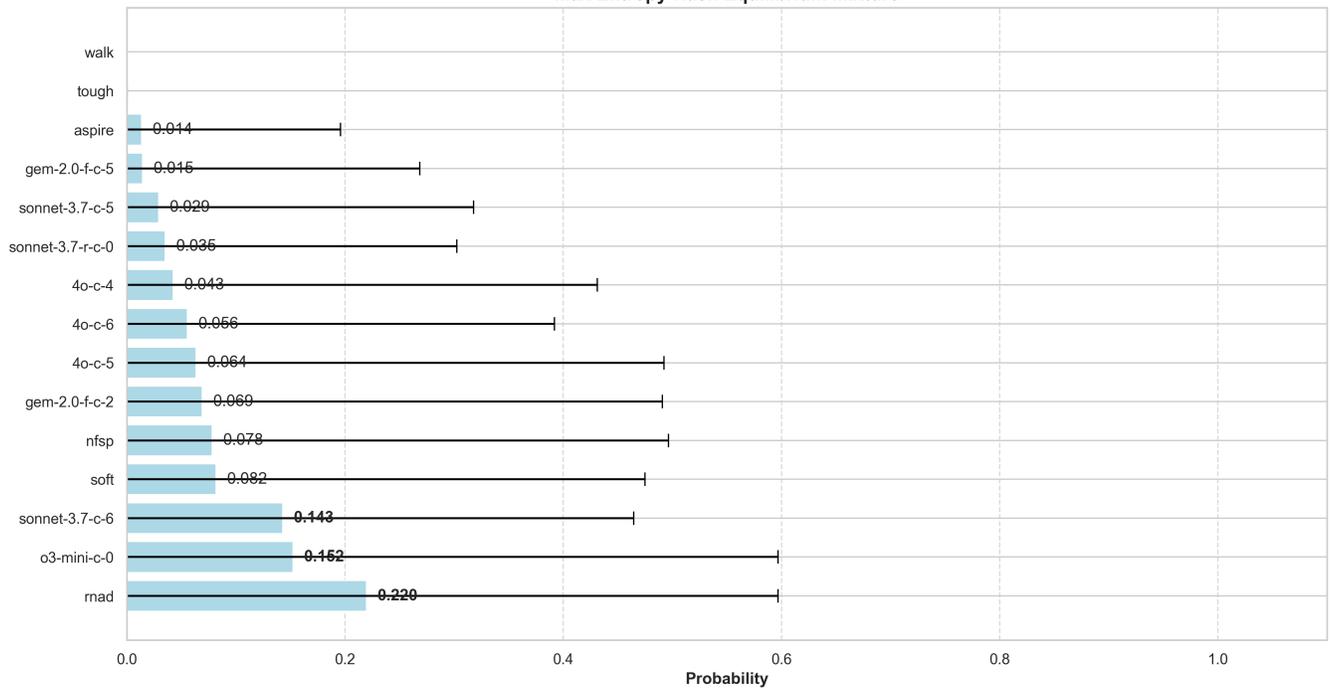


Figure 16: BG4, LLM+baselines

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

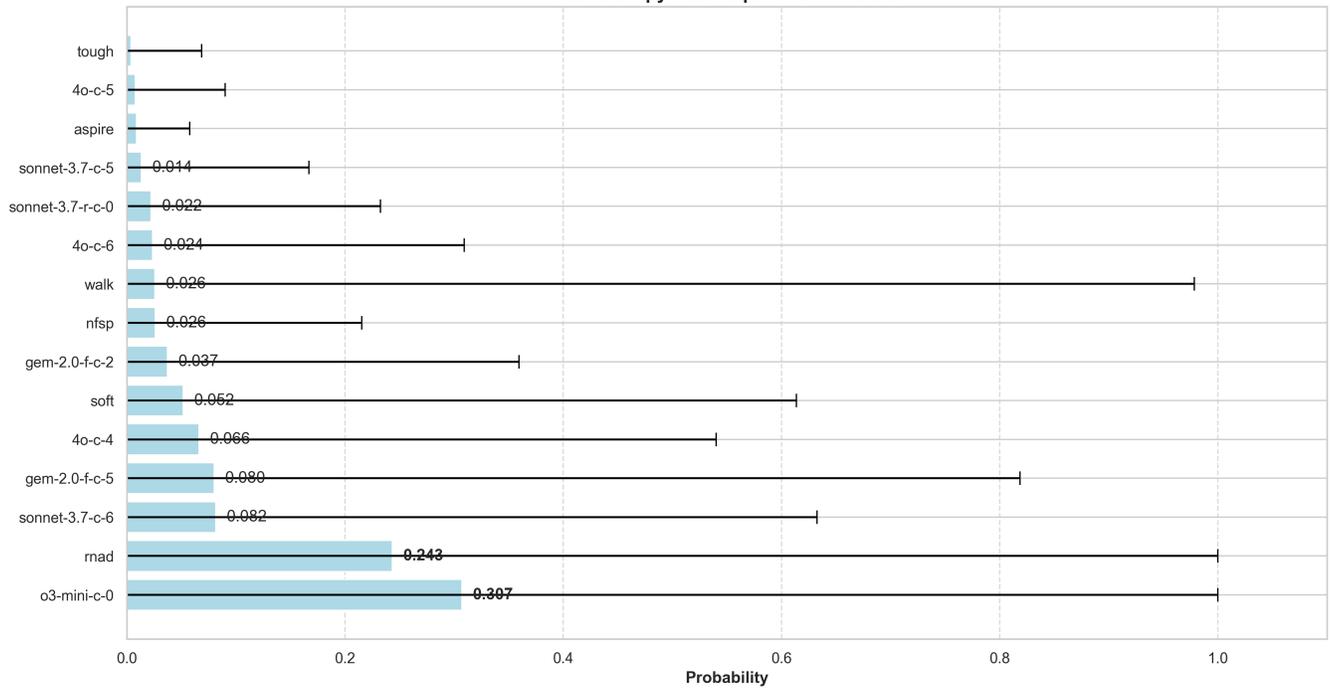


Figure 17: BG5, LLM+baselines

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

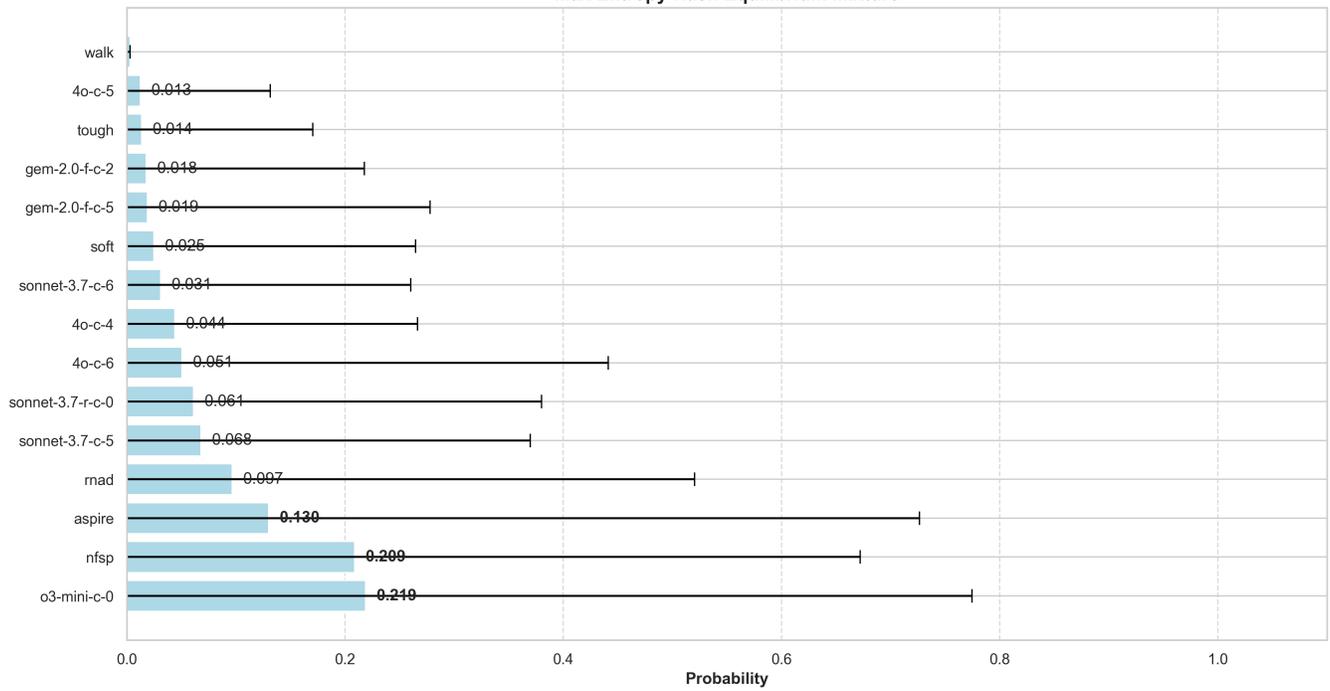


Figure 18: BG6, LLM+baselines

Average Nash Equilibrium Mixes (95% Confidence Intervals)

Max Entropy Nash Equilibrium Mixture

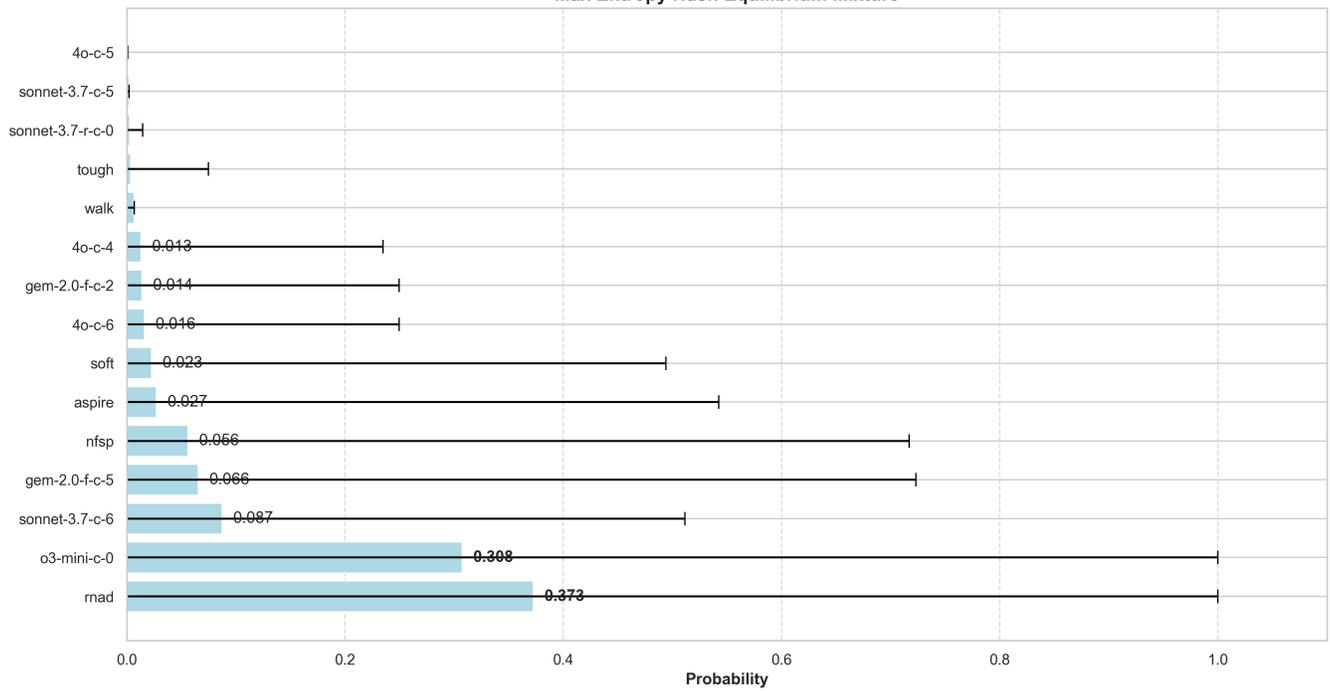


Figure 19: BGS0, LLM+baselines

F Disaggregated bootstrapping results for LLM-only and LLM+baselines

Table 2: Bootstrap means (\pm standard errors) of regret, normalized Nash, NW^+ , and utilitarian welfares, and EF1 frequencies, the last 4 expressed as percent points, for our LLM-only analysis on BG1-3. Lower values are better for regret and worse for the other metrics. The best value of each metric is in boldface.

| Agent | BG1 ($\gamma = 0.9, R = 3$) | | | | | BG2 ($\gamma = 0.98, R = 3$) | | | | |
|------------------|-----------------------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
| | MENE Regret | NW% | NW^+ % | UW % | EF1 % | MENE Regret | NW% | NW^+ % | UW % | EF1 % |
| sonnet-3.7-c-5 | 34.95 \pm 0.76 | 89.19 \pm 0.08 | 52.12 \pm 0.16 | 86.89 \pm 0.07 | 33.85 \pm 0.17 | 37.11 \pm 0.86 | 87.67 \pm 0.07 | 55.44 \pm 0.14 | 85.68 \pm 0.06 | 35.35 \pm 0.12 |
| sonnet-3.7-c-6 | 20.77 \pm 0.62 | 90.93 \pm 0.08 | 48.75 \pm 0.16 | 87.79 \pm 0.06 | 33.05 \pm 0.19 | 17.71 \pm 0.63 | 90.93 \pm 0.08 | 58.34 \pm 0.18 | 86.05 \pm 0.07 | 36.83 \pm 0.14 |
| sonnet-3.7-r-c-0 | 39.79 \pm 0.88 | 87.66 \pm 0.07 | 45.98 \pm 0.13 | 84.11 \pm 0.05 | 41.71 \pm 0.15 | 15.93 \pm 0.78 | 89.44 \pm 0.08 | 53.56 \pm 0.16 | 85.52 \pm 0.06 | 44.56 \pm 0.14 |
| gem-2.0-f-c-2 | 38.31 \pm 0.89 | 84.76 \pm 0.07 | 42.40 \pm 0.13 | 82.47 \pm 0.06 | 34.79 \pm 0.17 | 28.71 \pm 0.82 | 85.28 \pm 0.08 | 49.60 \pm 0.17 | 82.96 \pm 0.06 | 33.64 \pm 0.19 |
| gem-2.0-f-c-5 | 48.90 \pm 0.71 | 84.01 \pm 0.10 | 43.53 \pm 0.15 | 84.01 \pm 0.06 | 26.81 \pm 0.23 | 29.01 \pm 0.85 | 83.74 \pm 0.10 | 45.01 \pm 0.15 | 84.71 \pm 0.07 | 33.01 \pm 0.27 |
| 4o-c-4 | 13.27 \pm 0.56 | 91.97 \pm 0.08 | 49.14 \pm 0.13 | 86.62 \pm 0.06 | 37.88 \pm 0.20 | 13.27 \pm 0.56 | 90.42 \pm 0.08 | 52.26 \pm 0.18 | 85.72 \pm 0.07 | 41.80 \pm 0.17 |
| 4o-c-5 | 6.05 \pm 0.40 | 93.00 \pm 0.08 | 50.96 \pm 0.18 | 88.14 \pm 0.06 | 37.38 \pm 0.21 | 9.37 \pm 0.57 | 90.46 \pm 0.08 | 55.07 \pm 0.17 | 85.87 \pm 0.06 | 40.54 \pm 0.16 |
| 4o-c-6 | 20.49 \pm 0.67 | 92.83 \pm 0.07 | 54.34 \pm 0.17 | 87.66 \pm 0.05 | 38.86 \pm 0.26 | 18.49 \pm 0.74 | 89.33 \pm 0.06 | 59.28 \pm 0.18 | 85.38 \pm 0.05 | 38.38 \pm 0.15 |
| o3-mini-c-0 | 4.68 \pm 0.35 | 92.48 \pm 0.08 | 56.70 \pm .20 | 86.73 \pm 0.06 | 52.20 \pm 0.19 | 2.73 \pm 0.33 | 91.87 \pm 0.06 | 64.39 \pm 0.23 | 86.95 \pm 0.05 | 54.05 \pm 0.17 |

| Agent | BG3 ($\gamma = 0.98, R = 5$) | | | | |
|------------------|-----------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
| | MENE Regret | NW% | NW^+ % | UW % | EF1 % |
| sonnet-3.7-c-5 | 25.83 \pm 0.93 | 84.91 \pm 0.09 | 54.41 \pm .20 | 84.76 \pm 0.07 | 35.24 \pm 0.20 |
| sonnet-3.7-c-6 | 29.57 \pm 0.85 | 86.60 \pm 0.13 | 56.64 \pm 0.19 | 85.88 \pm 0.08 | 32.24 \pm 0.21 |
| sonnet-3.7-r-c-0 | 13.71 \pm 0.70 | 90.61 \pm 0.10 | 52.99 \pm 0.20 | 86.74 \pm 0.07 | 43.34 \pm 0.17 |
| gem-2.0-f-c-2 | 41.49 \pm 1.36 | 85.33 \pm 0.15 | 42.73 \pm 0.22 | 82.56 \pm 0.10 | 25.04 \pm 0.21 |
| gem-2.0-f-c-5 | 55.44 \pm 1.32 | 81.12 \pm 0.20 | 52.33 \pm 0.28 | 81.99 \pm 0.13 | 34.52 \pm 0.43 |
| 4o-c-4 | 25.59 \pm 0.83 | 88.79 \pm 0.07 | 57.44 \pm 0.21 | 83.81 \pm 0.06 | 41.53 \pm 0.18 |
| 4o-c-5 | 22.69 \pm 0.84 | 88.92 \pm 0.14 | 54.42 \pm 0.20 | 85.10 \pm 0.10 | 37.37 \pm 0.21 |
| 4o-c-6 | 7.69 \pm 0.56 | 88.66 \pm 0.08 | 54.84 \pm 0.18 | 84.64 \pm 0.06 | 38.54 \pm 0.20 |
| o3-mini-c-0 | 4.24 \pm 0.45 | 92.08 \pm 0.11 | 59.58 \pm 0.24 | 87.28 \pm 0.09 | 54.62 \pm 0.25 |

G NFSP training parameters for games in the BGS class

Table 6: NFSP training configuration for BG4 ($\gamma = 0.9, R = 3$)

| NFSP / training hyper-parameters | |
|---|--------------------|
| anticipatory parameter η | 0.1 |
| ϵ -greedy exploration (ϵ) | 0.06 |
| batch_size | 128 |
| policy_buffer_size | 100,000 |
| q_buffer_size | 20,000 |
| updates_per_loop | 3 |
| policy optimiser | Adam |
| Q-network optimiser | Adam |
| policy learning rate | 1×10^{-3} |
| Q-network learning rate | 1×10^{-3} |
| hidden layer size (both nets) | 256 |
| policy loss | Cross-Entropy |
| Q-network loss | Mean-Squared Error |
| Number of Game Trajectories | 10^8 |

Table 3: Bootstrap means (\pm standard errors) of regret, normalized Nash, NW^+ , and utilitarian welfares, and EF1 frequencies, the last 4 expressed as percent points, for our LLM-only analysis on BG4-6 and BGS0. Lower values are better for regret and worse for the other metrics. The best value of each metric is in boldface.

| Agent | BG4 ($\gamma = 0.9, R = 3$) | | | | | BG5 ($\gamma = 0.98, R = 3$) | | | | |
|------------------|-----------------------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|-----------------------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|
| | MENE Regret | NW% | NW^+ % | UW% | EF1% | MENE Regret | NW% | NW^+ % | UW% | EF1% |
| sonnet-3.7-c-5 | 26.24 \pm 0.61 | 79.66 \pm 0.08 | 48.56 \pm .21 | 87.14 \pm 0.06 | 27.59 \pm 0.19 | 28.77 \pm 0.73 | 78.50 \pm 0.14 | 57.05 \pm .28 | 86.70 \pm 0.08 | 31.05 \pm 0.28 |
| sonnet-3.7-c-6 | 18.34 \pm 0.57 | 79.25 \pm 0.07 | 57.27 \pm .22 | 86.82 \pm 0.07 | 29.45 \pm 0.16 | 22.64 \pm 0.83 | 77.90 \pm 0.14 | 51.62 \pm .17 | 86.66 \pm 0.09 | 33.49 \pm 0.35 |
| sonnet-3.7-r-c-0 | 15.39 \pm 0.58 | 82.26 \pm 0.07 | 55.92 \pm .21 | 88.70 \pm 0.05 | 38.38 \pm 0.23 | 17.42 \pm 0.61 | 82.51 \pm 0.11 | 57.52 \pm .21 | 89.69 \pm 0.08 | 40.44 \pm 0.31 |
| gemini-2.0-f-c-2 | 15.65 \pm 0.55 | 80.61 \pm 0.07 | 51.39 \pm .31 | 87.15 \pm 0.05 | 32.95 \pm 0.27 | 20.30 \pm 0.64 | 79.97 \pm 0.12 | 48.13 \pm .19 | 87.76 \pm 0.08 | 34.70 \pm 0.27 |
| gemini-2.0-f-c-5 | 27.89 \pm 0.55 | 78.34 \pm 0.11 | 51.05 \pm .18 | 86.84 \pm 0.07 | 28.36 \pm 0.27 | 18.17 \pm 0.58 | 75.53 \pm 0.21 | 44.38 \pm .18 | 86.34 \pm 0.12 | 23.63 \pm 0.24 |
| 4o-c-4 | 15.12 \pm 0.59 | 81.51 \pm 0.07 | 48.75 \pm .20 | 86.61 \pm 0.06 | 31.41 \pm 0.25 | 14.26 \pm 0.57 | 82.96 \pm 0.10 | 56.04 \pm .28 | 88.31 \pm 0.08 | 34.84 \pm 0.35 |
| 4o-c-5 | 16.88 \pm 0.52 | 81.98 \pm 0.11 | 53.56 \pm .27 | 87.83 \pm 0.09 | 31.35 \pm 0.27 | 27.37 \pm 0.66 | 79.25 \pm 0.13 | 52.62 \pm .17 | 86.36 \pm 0.08 | 32.13 \pm 0.28 |
| 4o-c-6 | 23.18 \pm 0.56 | 81.24 \pm 0.08 | 46.86 \pm .15 | 87.17 \pm 0.07 | 28.80 \pm 0.20 | 31.05 \pm 0.61 | 79.79 \pm 0.11 | 55.74 \pm .21 | 85.94 \pm 0.07 | 34.75 \pm 0.23 |
| o3-mini-c-0 | 0.50 \pm 0.10 | 83.59 \pm 0.07 | 63.21 \pm .30 | 87.45 \pm 0.06 | 51.57 \pm 0.21 | 0.94 \pm 0.16 | 85.45 \pm 0.10 | 72.04 \pm .49 | 90.01 \pm 0.08 | 57.34 \pm 0.33 |

| Agent | BG6 ($\gamma = 0.98, R = 5$) | | | | | BGS0 (uniform mixture of BG4-BG6) | | | | |
|------------------|-----------------------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|-----------------------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|
| | MENE Regret | NW% | NW^+ % | UW% | EF1% | MENE Regret | NW% | NW^+ % | UW% | EF1% |
| sonnet-3.7-c-5 | 10.48 \pm 0.50 | 78.29 \pm 0.11 | 54.87 \pm .21 | 86.83 \pm 0.07 | 30.03 \pm 0.27 | 28.64 \pm 0.50 | 79.85 \pm 0.13 | 53.49 \pm .14 | 86.94 \pm 0.06 | 34.20 \pm 0.28 |
| sonnet-3.7-c-6 | 21.31 \pm 0.58 | 76.30 \pm 0.11 | 55.90 \pm .20 | 85.47 \pm 0.06 | 27.95 \pm 0.25 | 26.09 \pm 0.54 | 78.13 \pm 0.12 | 54.16 \pm .13 | 86.67 \pm 0.06 | 34.15 \pm 0.29 |
| sonnet-3.7-r-c-0 | 16.17 \pm 0.65 | 80.64 \pm 0.10 | 54.32 \pm .22 | 86.94 \pm 0.06 | 38.35 \pm 0.26 | 22.87 \pm 0.50 | 82.34 \pm 0.09 | 55.71 \pm .15 | 88.51 \pm 0.05 | 44.18 \pm 0.33 |
| gemini-2.0-f-c-2 | 20.25 \pm 0.60 | 79.80 \pm 0.15 | 54.67 \pm .19 | 86.33 \pm 0.09 | 32.40 \pm 0.22 | 22.50 \pm 0.46 | 80.74 \pm 0.13 | 54.19 \pm .20 | 87.04 \pm 0.06 | 38.52 \pm 0.28 |
| gemini-2.0-f-c-5 | 21.80 \pm 0.70 | 73.78 \pm 0.11 | 54.65 \pm .15 | 83.77 \pm 0.07 | 27.38 \pm 0.26 | 22.48 \pm 0.54 | 77.63 \pm 0.18 | 49.46 \pm .12 | 86.24 \pm 0.07 | 31.81 \pm 0.31 |
| 4o-c-4 | 14.02 \pm 0.54 | 80.31 \pm 0.07 | 54.12 \pm .19 | 85.54 \pm 0.06 | 31.16 \pm 0.23 | 15.74 \pm 0.42 | 82.08 \pm 0.07 | 53.45 \pm .16 | 86.72 \pm 0.04 | 38.19 \pm 0.36 |
| 4o-c-5 | 16.35 \pm 0.55 | 79.26 \pm 0.11 | 54.21 \pm .17 | 86.24 \pm 0.06 | 31.33 \pm 0.30 | 22.09 \pm 0.40 | 81.08 \pm 0.12 | 56.09 \pm .17 | 87.20 \pm 0.06 | 38.04 \pm 0.33 |
| 4o-c-6 | 13.68 \pm 0.54 | 80.83 \pm 0.12 | 49.36 \pm .17 | 87.58 \pm 0.07 | 35.88 \pm 0.35 | 22.72 \pm 0.42 | 81.01 \pm 0.11 | 50.44 \pm .11 | 86.82 \pm 0.06 | 38.07 \pm 0.29 |
| o3-mini-c-0 | 1.50 \pm 0.19 | 85.04 \pm 0.11 | 60.44 \pm .25 | 89.08 \pm 0.10 | 51.98 \pm 0.25 | 0.12 \pm 0.04 | 86.34 \pm 0.11 | 68.19 \pm .33 | 89.76 \pm 0.07 | 57.40 \pm 0.30 |

Table 7: NFSP training configuration for BG5 ($\gamma = 0.98, R = 3$)

| NFSP / training hyper-parameters | |
|---|--------------------|
| anticipatory parameter η | 0.1 |
| ϵ -greedy exploration (ϵ) | 0.06 |
| batch_size | 128 |
| policy_buffer_size | 1,000,000 |
| q_buffer_size | 500,000 |
| updates_per_loop | 6 |
| policy optimiser | Adam |
| Q-network optimiser | Adam |
| policy learning rate | 1×10^{-5} |
| Q-network learning rate | 1×10^{-3} |
| hidden layer size (both nets) | 256 |
| policy loss | Cross-Entropy |
| Q-network loss | Mean-Squared Error |
| Number of Game Trajectories | 10^8 |

Table 4: Bootstrap means (\pm standard errors) of regret, normalized Nash, NW^+ , and utilitarian welfares, and EF1 frequencies, the last 4 expressed as percent points, for our LLM+baselines analysis on BG1–3. Lower values are better for regret and worse for the other metrics. The best value of each metric is in boldface.

| Agent | BG1 ($\gamma = 0.9, R = 3$) | | | | | BG2 ($\gamma = 0.98, R = 3$) | | | | |
|------------------|-----------------------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|-----------------------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|
| | MENE Regret | NW% | NW^+ % | UW% | EF1% | MENE Regret | NW% | NW^+ % | UW% | EF1% |
| sonnet-3.7-c-5 | 40.72 \pm 1.33 | 88.41 \pm 0.17 | 50.08 \pm .32 | 86.77 \pm 0.08 | 33.18 \pm 0.20 | 48.12 \pm 1.38 | 84.93 \pm 0.18 | 50.85 \pm .27 | 84.72 \pm 0.07 | 32.75 \pm 0.18 |
| sonnet-3.7-c-6 | 26.69 \pm 1.33 | 89.59 \pm 0.20 | 47.00 \pm .29 | 87.31 \pm 0.08 | 32.21 \pm 0.20 | 27.34 \pm 1.19 | 86.22 \pm 0.19 | 56.07 \pm .24 | 85.47 \pm 0.08 | 34.98 \pm 0.20 |
| sonnet-3.7-r-c-0 | 36.73 \pm 0.85 | 87.19 \pm 0.14 | 44.31 \pm .27 | 84.04 \pm 0.07 | 41.61 \pm 0.19 | 14.21 \pm 0.72 | 87.30 \pm 0.15 | 53.23 \pm .27 | 84.47 \pm 0.08 | 42.97 \pm 0.22 |
| gemini-2.0-f-c-2 | 37.69 \pm 0.88 | 84.94 \pm 0.08 | 40.46 \pm .26 | 82.84 \pm 0.06 | 34.39 \pm 0.20 | 30.64 \pm 0.90 | 83.45 \pm 0.14 | 46.20 \pm .23 | 82.17 \pm 0.08 | 30.52 \pm 0.24 |
| gemini-2.0-f-c-5 | 50.76 \pm 0.90 | 84.16 \pm 0.13 | 41.57 \pm .28 | 83.74 \pm 0.08 | 25.72 \pm 0.24 | 32.40 \pm 0.92 | 81.91 \pm 0.13 | 41.21 \pm .24 | 83.85 \pm 0.08 | 29.40 \pm 0.29 |
| 4o-c-4 | 14.65 \pm 0.55 | 91.87 \pm 0.09 | 47.11 \pm .29 | 86.59 \pm 0.07 | 36.27 \pm 0.22 | 19.80 \pm 0.76 | 88.43 \pm 0.13 | 51.06 \pm .25 | 84.44 \pm 0.09 | 38.19 \pm 0.26 |
| 4o-c-5 | 5.91 \pm 0.43 | 92.48 \pm 0.11 | 48.61 \pm .32 | 87.66 \pm 0.08 | 36.51 \pm 0.24 | 9.32 \pm 0.62 | 89.26 \pm 0.10 | 51.94 \pm .24 | 85.42 \pm 0.07 | 37.25 \pm 0.22 |
| 4o-c-6 | 17.77 \pm 0.74 | 92.78 \pm 0.10 | 52.07 \pm .33 | 87.72 \pm 0.08 | 37.37 \pm 0.26 | 15.82 \pm 0.72 | 88.95 \pm 0.10 | 54.73 \pm .30 | 85.35 \pm 0.08 | 34.67 \pm 0.23 |
| o3-mini-c-0 | 4.37 \pm 0.34 | 92.25 \pm 0.10 | 54.81 \pm .36 | 86.69 \pm 0.07 | 50.84 \pm 0.26 | 4.07 \pm 0.43 | 90.47 \pm 0.12 | 60.38 \pm .34 | 85.68 \pm 0.10 | 50.44 \pm 0.27 |
| soft | 106.52 \pm 1.81 | 79.95 \pm 0.28 | 36.00 \pm .61 | 77.13 \pm 0.20 | 21.15 \pm 0.43 | 46.91 \pm 1.72 | 80.23 \pm 0.28 | 40.81 \pm .53 | 79.99 \pm 0.16 | 22.64 \pm 0.39 |
| tough | 108.13 \pm 2.23 | 76.11 \pm 0.28 | 1.62 \pm .05 | 79.18 \pm 0.18 | 0.00 \pm 0.00 | 67.82 \pm 2.05 | 63.60 \pm 0.19 | 1.85 \pm .03 | 73.46 \pm 0.14 | 0.00 \pm 0.00 |
| walk | 79.65 \pm 0.73 | 72.01 \pm 0.00 | 0.00 \pm .00 | 72.08 \pm 0.00 | 0.00 \pm 0.00 | 94.83 \pm 0.71 | 72.00 \pm 0.00 | 0.00 \pm .00 | 72.08 \pm 0.00 | 0.00 \pm 0.00 |

| Agent | BG3 ($\gamma = 0.98, R = 5$) | | | | |
|------------------|-----------------------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|
| | MENE Regret | NW% | NW^+ % | UW% | EF1% |
| sonnet-3.7-c-5 | 52.47 \pm 1.90 | 83.21 \pm 0.25 | 54.80 \pm .46 | 83.87 \pm 0.10 | 35.47 \pm 0.25 |
| sonnet-3.7-c-6 | 61.69 \pm 2.19 | 85.66 \pm 0.26 | 54.87 \pm .36 | 85.13 \pm 0.12 | 34.79 \pm 0.28 |
| sonnet-3.7-r-c-0 | 26.38 \pm 1.35 | 88.57 \pm 0.23 | 47.39 \pm .31 | 85.36 \pm 0.13 | 42.79 \pm 0.28 |
| gemini-2.0-f-c-2 | 56.86 \pm 1.69 | 84.28 \pm 0.18 | 41.15 \pm .30 | 80.90 \pm 0.13 | 27.81 \pm 0.23 |
| gemini-2.0-f-c-5 | 71.36 \pm 2.04 | 79.17 \pm 0.30 | 47.28 \pm .39 | 81.41 \pm 0.16 | 36.03 \pm 0.50 |
| 4o-c-4 | 37.12 \pm 1.25 | 88.16 \pm 0.14 | 55.82 \pm .37 | 83.49 \pm 0.11 | 38.22 \pm 0.29 |
| 4o-c-5 | 29.71 \pm 1.43 | 87.87 \pm 0.21 | 52.08 \pm .35 | 84.18 \pm 0.16 | 36.97 \pm 0.29 |
| 4o-c-6 | 17.88 \pm 1.06 | 89.05 \pm 0.21 | 50.43 \pm .36 | 85.73 \pm 0.18 | 38.11 \pm 0.27 |
| o3-mini-c-0 | 3.25 \pm 0.46 | 92.82 \pm 0.16 | 57.90 \pm .40 | 87.71 \pm 0.13 | 54.68 \pm 0.37 |
| soft | 33.95 \pm 1.61 | 88.18 \pm 0.40 | 45.62 \pm .69 | 85.69 \pm 0.25 | 32.00 \pm 0.42 |
| tough | 89.61 \pm 2.37 | 60.63 \pm 0.28 | 4.94 \pm .08 | 71.98 \pm 0.15 | 0.32 \pm 0.03 |
| walk | 104.13 \pm 1.10 | 72.01 \pm 0.00 | 0.00 \pm .00 | 72.08 \pm 0.00 | 0.00 \pm 0.00 |

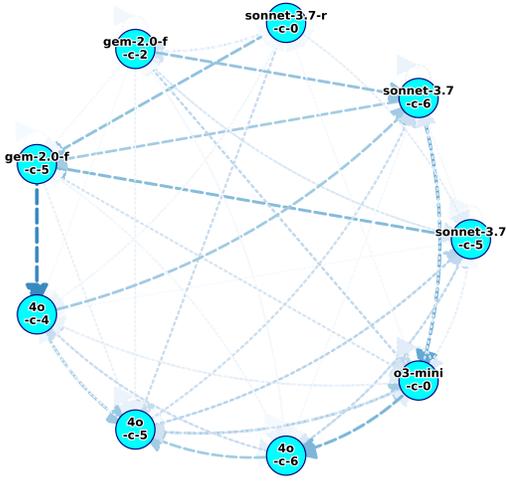
Table 8: NFSP training configuration for BG6 ($\gamma = 0.98, R = 3$)

| NFSP / training hyper-parameters | |
|---|--------------------|
| anticipatory parameter η | 0.1 |
| ϵ -greedy exploration (ϵ) | 0.06 |
| batch_size | 128 |
| policy_buffer_size | 2,000,000 |
| q_buffer_size | 1,000,000 |
| updates_per_loop | 10 |
| policy optimiser | Adam |
| Q-network optimiser | Adam |
| policy learning rate | 1×10^{-5} |
| Q-network learning rate | 1×10^{-3} |
| hidden layer size (both nets) | 256 |
| policy loss | Cross-Entropy |
| Q-network loss | Mean-Squared Error |
| Number of Game Trajectories | 10^8 |

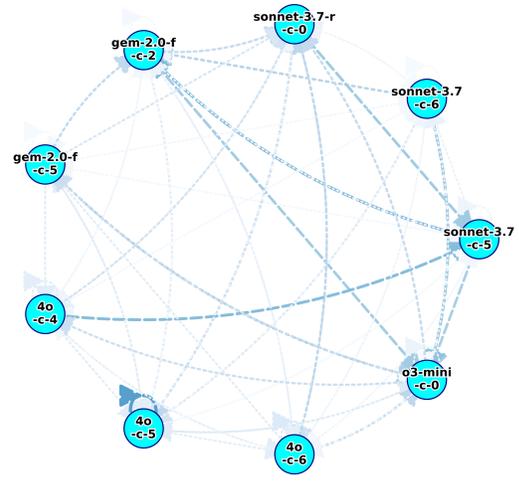
Table 5: Bootstrap means (\pm standard errors) of regret, normalized Nash, NW^+ , and utilitarian welfares, and EF1 frequencies, the last 4 expressed as percent points, for our LLM+baselines analysis on BG4–6. Lower values are better for regret and worse for the other metrics. The best value of each metric is in boldface.

| Agent | BG4 ($\gamma = 0.9, R = 3$) | | | | | BG5 ($\gamma = 0.98, R = 3$) | | | | |
|------------------|--------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|--------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|
| | MENE Regret | NW% | NW^+ % | UW % | EF1 % | MENE Regret | NW% | NW^+ % | UW % | EF1 % |
| sonnet-3.7-c-5 | 17.19 \pm .44 | 78.79 \pm .07 | 49.43 \pm .16 | 86.97 \pm .05 | 29.61 \pm .15 | 43.27 \pm .81 | 76.54 \pm .20 | 55.89 \pm .38 | 85.22 \pm .10 | 29.16 \pm .28 |
| sonnet-3.7-c-6 | 6.82 \pm .36 | 77.77 \pm .07 | 54.04 \pm .20 | 86.41 \pm .06 | 30.87 \pm .16 | 27.12 \pm .80 | 78.08 \pm .15 | 54.50 \pm .38 | 86.32 \pm .10 | 32.32 \pm .35 |
| sonnet-3.7-r-c-0 | 18.08 \pm .54 | 80.13 \pm .08 | 50.23 \pm .18 | 86.68 \pm .06 | 41.43 \pm .18 | 21.03 \pm .58 | 80.62 \pm .12 | 58.44 \pm .41 | 87.55 \pm .11 | 39.21 \pm .33 |
| gemini-2.0-f-c-2 | 17.86 \pm .46 | 78.05 \pm .10 | 47.94 \pm .20 | 85.81 \pm .07 | 35.45 \pm .21 | 18.93 \pm .57 | 80.06 \pm .15 | 50.62 \pm .35 | 87.67 \pm .10 | 34.13 \pm .28 |
| gemini-2.0-f-c-5 | 24.28 \pm .39 | 76.94 \pm .20 | 51.54 \pm .26 | 86.14 \pm .06 | 29.38 \pm .19 | 22.96 \pm .59 | 75.94 \pm .22 | 44.09 \pm .29 | 86.52 \pm .12 | 25.00 \pm .26 |
| 4o-c-4 | 17.60 \pm .47 | 79.92 \pm .07 | 50.00 \pm .14 | 85.90 \pm .06 | 33.74 \pm .21 | 16.40 \pm .50 | 82.01 \pm .13 | 55.80 \pm .39 | 87.45 \pm .11 | 34.41 \pm .36 |
| 4o-c-5 | 14.74 \pm .40 | 79.97 \pm .09 | 52.35 \pm .19 | 86.67 \pm .08 | 34.65 \pm .22 | 26.42 \pm .58 | 75.94 \pm .22 | 49.34 \pm .31 | 86.12 \pm .10 | 31.50 \pm .29 |
| 4o-c-6 | 11.22 \pm .36 | 80.81 \pm .09 | 45.64 \pm .12 | 87.96 \pm .08 | 30.27 \pm .16 | 23.91 \pm .55 | 80.00 \pm .14 | 51.20 \pm .33 | 86.47 \pm .11 | 31.96 \pm .26 |
| o3-mini-c-0 | 1.60 \pm .20 | 82.62\pm.07 | 56.79\pm.21 | 87.56 \pm .05 | 50.76 \pm .17 | 2.36\pm.23 | 85.65\pm.14 | 66.83\pm.55 | 90.05\pm.11 | 54.79\pm.41 |
| nfsp | 13.58 \pm .51 | 77.91 \pm .08 | 44.73 \pm .16 | 84.00 \pm .06 | 40.57 \pm .23 | 17.90 \pm .56 | 82.62 \pm .22 | 65.05 \pm .36 | 89.59 \pm .14 | 52.32 \pm .57 |
| rnad | 1.17\pm.16 | 80.49 \pm .07 | 53.70 \pm .20 | 87.90 \pm .07 | 53.72\pm.20 | 4.80 \pm .31 | 85.03 \pm .17 | 56.39 \pm .37 | 90.15 \pm .13 | 48.27 \pm .41 |
| aspire | 25.83 \pm .50 | 81.45 \pm .08 | 43.72 \pm .11 | 88.32\pm.07 | 30.73 \pm .09 | 24.61 \pm .47 | 84.03 \pm .13 | 62.49 \pm .36 | 88.29 \pm .10 | 46.53 \pm .37 |
| soft | 21.53 \pm .53 | 70.02 \pm .08 | 42.98 \pm .19 | 80.12 \pm .05 | 24.75 \pm .11 | 36.69 \pm .69 | 70.37 \pm .12 | 46.43 \pm .35 | 80.47 \pm .06 | 27.59 \pm .26 |
| tough | 44.77 \pm .53 | 63.50 \pm .16 | 9.51 \pm .05 | 79.09 \pm .06 | 1.55 \pm .05 | 49.97 \pm .83 | 64.34 \pm .18 | 6.69 \pm .08 | 77.45 \pm .07 | .07 \pm .41 |
| walk | 37.08 \pm .30 | 67.97 \pm 0.00 | 0.00 \pm 0.00 | 75.03 \pm 0.00 | 0.00 \pm 0.00 | 60.95 \pm .58 | 67.98 \pm 0.00 | 0.00 \pm 0.00 | 75.04 \pm 0.00 | 0.00 \pm 0.00 |

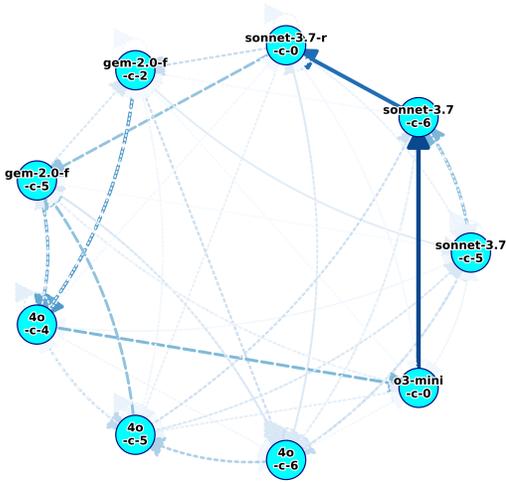
| Agent | BG6 ($\gamma = 0.98, R = 5$) | | | | |
|------------------|--------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|
| | MENE Regret | NW% | NW^+ % | UW % | EF1 % |
| sonnet-3.7-c-5 | 8.52 \pm .36 | 78.05 \pm .09 | 59.34 \pm .26 | 86.78 \pm .06 | 31.17 \pm .20 |
| sonnet-3.7-c-6 | 13.72 \pm .42 | 76.80 \pm .09 | 59.26 \pm .26 | 85.98 \pm .05 | 29.68 \pm .20 |
| sonnet-3.7-r-c-0 | 9.30 \pm .36 | 80.02 \pm .09 | 54.76 \pm .24 | 86.87 \pm .07 | 37.95 \pm .20 |
| gemini-2.0-f-c-2 | 17.01 \pm .41 | 79.57 \pm .09 | 51.43 \pm .19 | 86.52 \pm .07 | 32.21 \pm .17 |
| gemini-2.0-f-c-5 | 26.41 \pm .55 | 74.77 \pm .11 | 56.50 \pm .19 | 84.34 \pm .06 | 28.97 \pm .19 |
| 4o-c-4 | 7.42 \pm .32 | 81.23 \pm .08 | 58.44 \pm .25 | 86.56 \pm .06 | 34.81 \pm .25 |
| 4o-c-5 | 12.66 \pm .36 | 79.54 \pm .08 | 55.37 \pm .21 | 86.46 \pm .60 | 34.18 \pm .22 |
| 4o-c-6 | 13.90 \pm .39 | 80.10 \pm .08 | 49.63 \pm .17 | 86.60 \pm .60 | 33.33 \pm .24 |
| o3-mini-c-0 | 2.13\pm.23 | 84.05\pm.08 | 61.05 \pm .24 | 87.86 \pm .07 | 53.86 \pm .25 |
| nfsp | 2.85 \pm .25 | 82.33 \pm .08 | 65.79 \pm .36 | 87.42\pm.07 | 55.83\pm.30 |
| rnad | 15.42 \pm .52 | 81.25 \pm .12 | 51.75 \pm .20 | 86.50 \pm .08 | 42.94 \pm .27 |
| aspire | 11.26 \pm .41 | 83.85 \pm .08 | 67.23\pm.20 | 88.35 \pm .06 | 46.36 \pm .19 |
| soft | 20.52 \pm .50 | 69.32 \pm .11 | 46.27 \pm .19 | 80.13 \pm .04 | 25.18 \pm .15 |
| tough | 36.58 \pm .66 | 64.51 \pm .13 | 6.11 \pm .06 | 77.21 \pm .05 | 0.18 \pm .03 |
| walk | 46.55 \pm .33 | 67.97 \pm 0.00 | 0.00 \pm 0.00 | 75.04 \pm 0.04 | 0.00 \pm 0.00 |



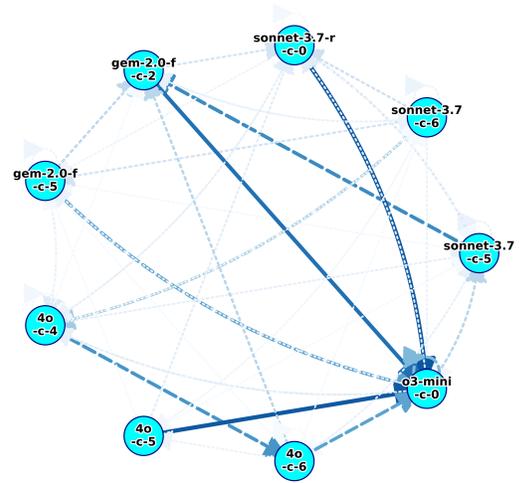
(a) BG1 ($\gamma = 0.90, R = 3$)



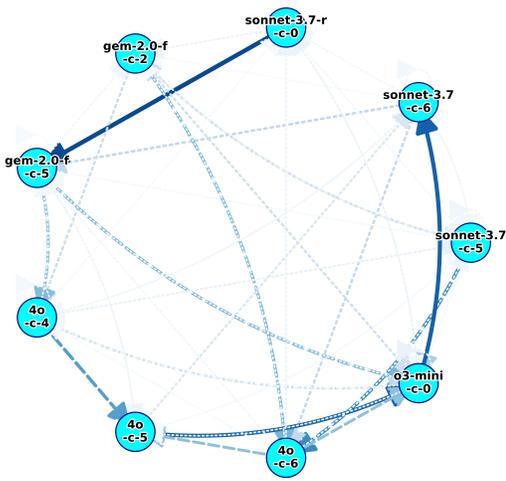
(b) BG4 ($\gamma = 0.90, R = 3$)



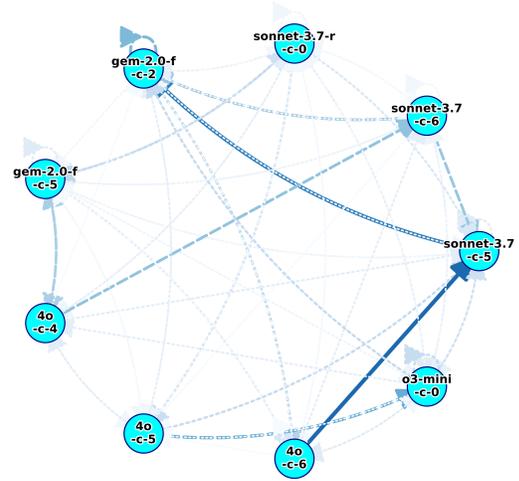
(c) BG2 ($\gamma = 0.98, R = 3$)



(d) BG5 ($\gamma = 0.98, R = 3$)

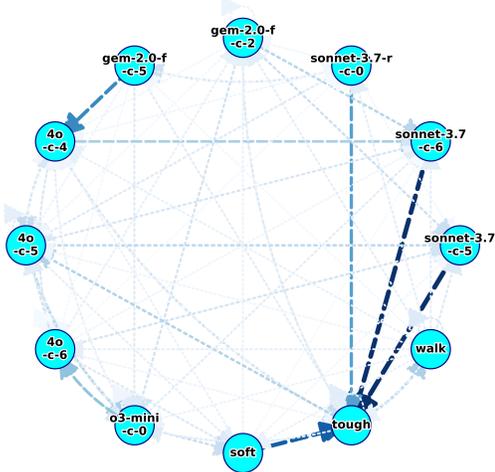


(e) BG3 ($\gamma = 0.98, R = 5$)

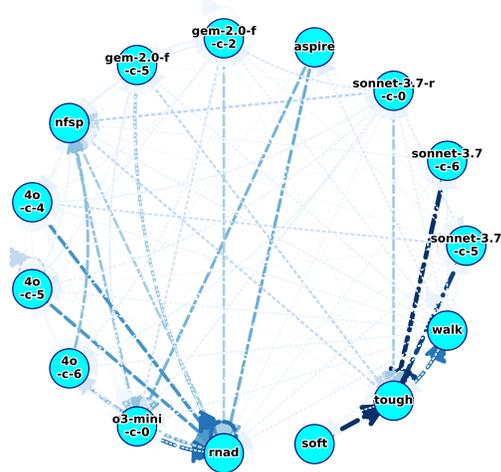


(f) BG6 ($\gamma = 0.98, R = 5$)

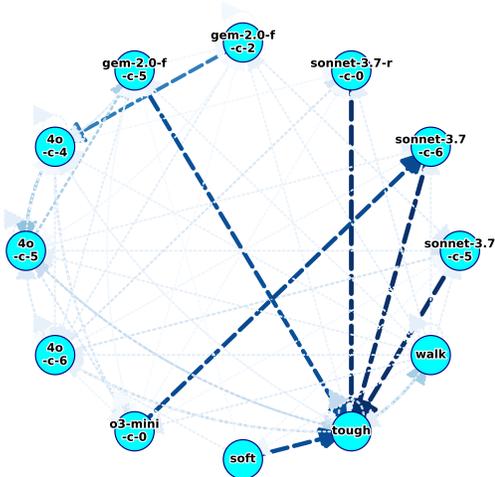
Figure 20: LLM-only empirical best-response graphs for BG1-6.



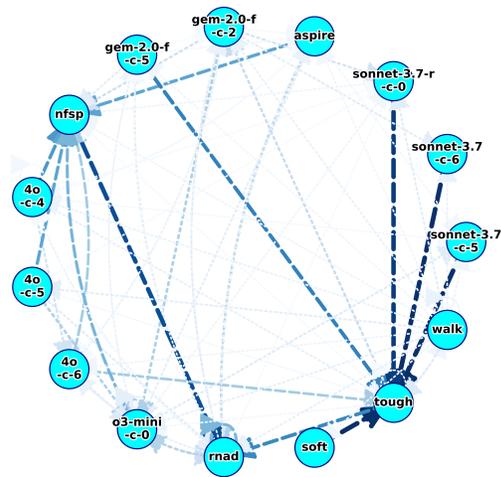
(a) BG1 ($\gamma = 0.90, R = 3$)



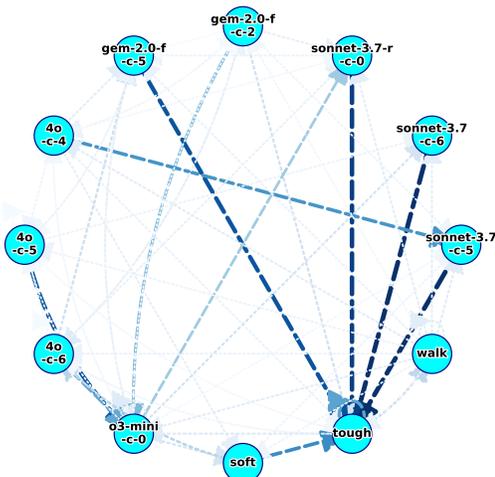
(b) BG4 ($\gamma = 0.90, R = 3$)



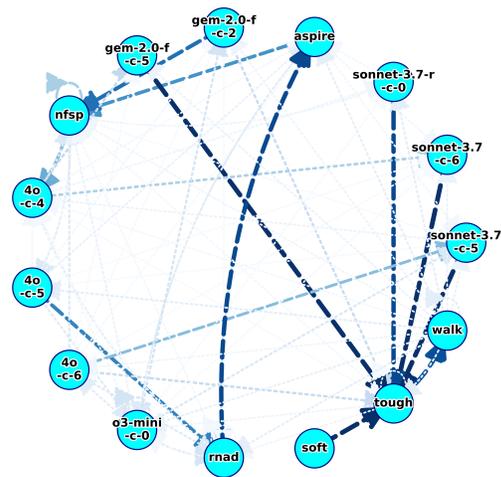
(c) BG2 ($\gamma = 0.98, R = 3$)



(d) BG5 ($\gamma = 0.98, R = 3$)



(e) BG3 ($\gamma = 0.98, R = 5$)



(f) BG6 ($\gamma = 0.98, R = 5$)

Figure 21: LLM+baselines empirical best-response graphs for BG1-6.

H R-NaD training parameters for games in the BGS class

Table 9: R-NaD training configuration for BG4 ($\gamma = 0.9$, $R = 3$)

| R-NaD training hyper-parameters | |
|--|--------------------------|
| trajectory_max | 6 |
| policy network hidden layers | [256, 256] |
| batch_size | 256 |
| learning_rate | 5×10^{-5} |
| clip_gradient | 1.0 |
| target_network_avg | 0.01 |
| training_steps | 1,000,000 |
| Adam ($\beta_1, \beta_2, \epsilon$) | (0.9, 0.999, 10^{-8}) |
| entropy_schedule_repeats | [1, 2, 1] |
| entropy_schedule_size | [5,000, 10,000, 20,000] |
| η (reward transform) | 0.2 |
| NeuRD β | 2.0 |
| NeuRD clip | 150 |
| V-trace c | 1.0 |
| fine-tune from_learner_steps | 800,000 |
| policy_threshold | 0.03 |
| policy_discretization | 32 |
| seed | 42 |

Table 10: R-NaD training configuration for BG5 ($\gamma = 0.98$, $R = 3$)

| R-NaD training hyper-parameters | |
|--|--------------------------|
| trajectory_max | 6 |
| policy network hidden layers | [256, 256] |
| batch_size | 256 |
| learning_rate | 5×10^{-5} |
| clip_gradient | 1.0 |
| target_network_avg | 0.01 |
| training_steps | 1,000,000 |
| Adam ($\beta_1, \beta_2, \epsilon$) | (0.9, 0.999, 10^{-8}) |
| entropy_schedule_repeats | [1, 2, 1] |
| entropy_schedule_size | [5,000, 10,000, 20,000] |
| η (reward transform) | 0.2 |
| NeuRD β | 2.0 |
| NeuRD clip | 150 |
| V-trace c | 1.0 |
| fine-tune from_learner_steps | 800,000 |
| policy_threshold | 0.03 |
| policy_discretization | 32 |
| seed | 42 |

Table 11: R-NaD training configuration for BG6 ($\gamma = 0.98$, $R = 5$)

| R-NaD training hyper-parameters | |
|--|--------------------------|
| trajectory_max | 10 |
| policy network hidden layers | [256, 256] |
| batch_size | 256 |
| learning_rate | 5×10^{-5} |
| clip_gradient | 1.0 |
| target_network_avg | 0.01 |
| training_steps | 1,000,000 |
| Adam ($\beta_1, \beta_2, \epsilon$) | (0.9, 0.999, 10^{-8}) |
| entropy_schedule_repeats | [1, 2, 1] |
| entropy_schedule_size | [5,000, 10,000, 20,000] |
| η (reward transform) | 0.2 |
| NeuRD β | 2.0 |
| NeuRD clip | 150 |
| V-trace c | 1.0 |
| fine-tune from_learner_steps | 800,000 |
| policy_threshold | 0.03 |
| policy_discretization | 32 |
| seed | 42 |

I Acceptance % vs. Tough

| Agent | Turn level acceptance rate | | Game level acceptance rate | |
|------------------|----------------------------|-------------|----------------------------|-------------|
| | vs tough | vs others | vs tough | vs others |
| sonnet-3.7-c-5 | 49.7% | 12.3% | 78.9% | 19.5% |
| sonnet-3.7-c-6 | 51.6% | 13.8% | 85.5% | 21.3% |
| sonnet-3.7-r-c-0 | 26.4% | 10.2% | 61.1% | 15.8% |
| gem-2.0-f-c-2 | 6.4% | 5.7% | 15.0% | 8.9% |
| gem-2.0-f-c-5 | 15.3% | 9.1% | 33.1% | 15.5% |
| 4o-c-4 | 1.9% | 1.8% | 5.5% | 2.7% |
| 4o-c-5 | 1.5% | 3.4% | 3.6% | 5.1% |
| 4o-c-6 | 2.6% | 3.0% | 6.0% | 4.5% |
| o3-mini-c-0 | 0.3% | 0.2% | 0.6% | 0.3% |
| soft | 100.0% | 100.0% | 100.0% | 100.0% |
| tough | 0.0% | 0.0% | 0.0% | 0.0% |
| walk | 0.0% | 0.0% | 0.0% | 0.0% |

Table 12: Acceptance of deals with offer value below the BATNA, reported at the turn level and the game level, and split by opponent type. Columns “vs tough” consider only negotiations against tough; “vs others” aggregates all other opponents.

This table reports mistaken acceptances: the frequency with which each agent accepts an offer whose value is strictly below its own BATNA. We report this error rate at two granularities. Turn-level acceptance rate measures, across all decision points where the agent had the option to accept, the fraction of turns on which it accepts a below-BATNA offer. Game-level acceptance rate measures, across complete negotiations, the fraction of games in which the agent ever makes such an acceptance (hence these rates are typically higher than turn-level rates). Results are split by opponent type: vs tough restricts to negotiations against the rigid tough heuristic, while “vs others” aggregates all remaining opponents.

J Opposition Levels

To better understand bargaining dynamics within our game, we compute the opposition levels for both BGS and BGL in Figs. 22 and 23. Opposition level captures the degree to which the two agents can simultaneously achieve high utility under the same terminal outcome. Following the definition used in NegMAS, we measure this by considering the achievable utility points in normalized utility space and computing the distance of the closest feasible point to the ideal point $(1, 1)$.

Concretely, let \mathcal{O} be the set of feasible terminal outcomes (allocations), including the disagreement/BATNA as a terminal outcome, and let $u_i(o)$ denote player i 's utility for the outcome o . Let $U_i^{\max} := \max_{o \in \mathcal{O}} u_i(o)$ be the maximum attainable utility for player i over \mathcal{O} . For each outcome o , we compute the normalized utilities $\tilde{u}_i(o) = u_i(o)/U_i^{\max}$ and define its distance to the ideal point as

$$d(o) = \sqrt{\sum_{i \in \{1,2\}} (1 - \tilde{u}_i(o))^2} = \sqrt{\left(1 - \frac{u_1(o)}{U_1^{\max}}\right)^2 + \left(1 - \frac{u_2(o)}{U_2^{\max}}\right)^2}.$$

The opposition level of the instance is then

$$\text{Opp} = \min_{o \in \mathcal{O}} d(o).$$

Intuitively, $\text{Opp} = 0$ when there exists an outcome that simultaneously achieves both players' maximum utilities, while larger values indicate stronger opposition because even the closest feasible point remains far from the ideal $(1, 1)$ in normalized utility space.

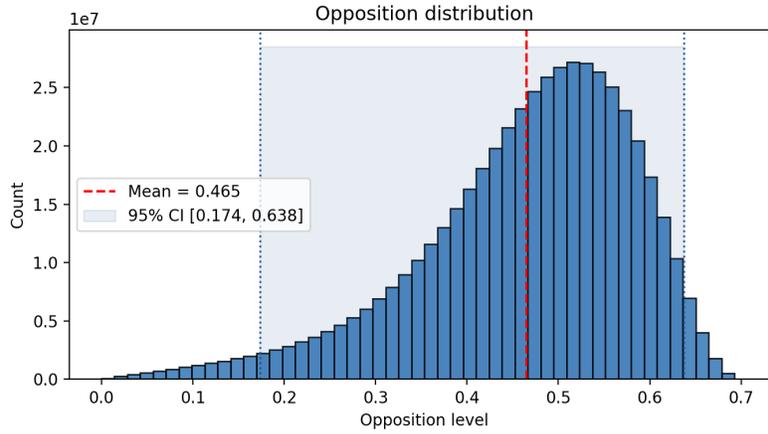


Figure 22: Opposition levels for the BGL setting of games.

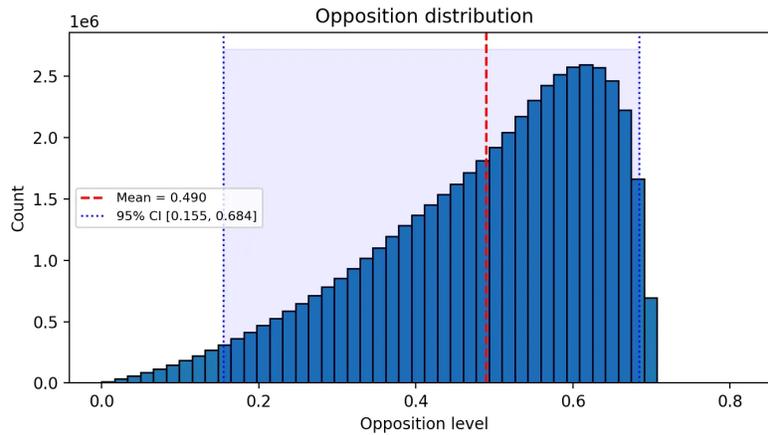


Figure 23: Opposition levels for the BGS setting of games.

K Provider Level Analysis

Table 13: Equilibrium weighted mean value given (mean \pm SE, $n = 1000$).

| Family | Mean \pm SE |
|-----------|------------------|
| Anthropic | 228.3 \pm 0.76 |
| Gemini | 213.2 \pm 1.94 |
| OpenAI | 222.4 \pm 0.87 |

Table 14: Empirical frequency with which the row provider is a best response to the column provider.

| | Anthropic | Gemini | OpenAI |
|-----------|-----------|--------|--------|
| Anthropic | 0.04 | 0.06 | 0.44 |
| Gemini | 0.73 | 0.00 | 0.00 |
| OpenAI | 0.23 | 0.94 | 0.56 |

The top table 13 gives the amount each model provider gives in equilibrium in the provider level meta-game. The table below 14 gives the empirical frequencies each provider is the best response to another.

K.1 Description of Provider Level Meta-Game

To construct the provider-level meta-game, we treat each LLM provider as a player. Each provider’s strategy set consists of the LLM-based meta-strategies in our library that are supplied by that provider.

We first group strategies by provider. For each ordered pair of providers (A, B) , we form the induced two-player subgame between A ’s strategies \mathcal{S}_A and B ’s strategies \mathcal{S}_B . Let $R \in \mathbb{R}^{|\mathcal{S}_A| \times |\mathcal{S}_B|}$ denote the corresponding payoff matrix for the row player, where each entry R_{ij} is the expected utility achieved by A ’s i -th strategy when negotiating against B ’s j -th strategy (estimated from cross-play).

We solve this subgame using the Lemke–Howson algorithm to obtain equilibrium mixtures $\sigma \in \Delta(\mathcal{S}_A)$ and $\tau \in \Delta(\mathcal{S}_B)$. The row provider’s equilibrium value for this matchup,

$$v(A, B) = \sigma^\top R \tau,$$

defines the (A, B) entry of the provider-level payoff matrix V . Repeating this procedure for all ordered pairs yields a square provider-level game whose size is the number of providers. We then analyze this provider-level game using the same best response graph pipeline as at the agent level, interpreting directed edges as empirical best responses between providers.

L Aspiration

Set up

- Rounds $t \in \{1, \dots, N\}$ with $N \in \{3, 5\}$.
- Discount per round γ^{t-1} , $\gamma \in \{0.9, 0.98\}$.
- x denotes the opponent's offer to us; y denotes our offer to the opponent.
- $\max_{\text{aspiration}} \in [0, 1]$ is a ceiling fraction. We have it set to 0.85 by default (NegMAS uses 1.0). It never becomes the opening target directly, instead we blend it with the sampled BATNA in every game (details below). This way our opening target scales with the BATNA.
- θ_t is the raw-utility target produced by the concession schedule. Given BATNA B and total value U_{max} we compute

$$\theta_1 = B + s(\gamma) \cdot \max_{\text{aspiration}} \cdot (U_{\text{max}} - B)$$

where the scale $s(\gamma)$ is a piecewise constant depending on the discount factor:

$$s(\gamma) = \begin{cases} 0.15 & \text{if } \gamma \leq 0.92 \\ 0.25 & \text{if } 0.92 < \gamma \leq 0.96 \\ 0.40 & \text{if } \gamma > 0.96 \end{cases}$$

This keeps θ_1 in $[B, U_{\text{max}})$, with patient games (high γ) starting at higher aspiration levels.

- $\epsilon = 0.05$ is the acceptance slack: at round t we shrink the target by 5%,

$$\theta_t^{\text{eff}} = \max(\theta_t - \epsilon, B),$$

so we never demand less than the BATNA.

Concession Schedules

After the dynamic opening value is set, later rounds follow a NegMAS curve that slides monotonically back to the BATNA. The exponent p is chosen from NegMAS presets based on γ , coupled with the initial scale $s(\gamma)$:

| γ | s | Profile | p (poly) |
|----------------|------|----------|------------|
| ≤ 0.92 | 0.15 | conceder | 0.25 |
| $(0.92, 0.96]$ | 0.25 | linear | 1.0 |
| > 0.96 | 0.40 | boulware | 4.0 |

This coupling ensures aggressive early concession when discounting is steep, and patient holding-out when future payoffs are nearly undiscounted. Let $\theta_N = B$ and $\Delta = \theta_1 - \theta_N$.

Polynomial family

$$\theta_t = \theta_N + \Delta \left(1 - \left(\frac{t-1}{N-1} \right)^p \right), \quad p \in \{4.0 \text{ (Boulware)}, 1.0 \text{ (linear)}, 0.25 \text{ (conceder)}\}.$$

Exponential family

$$\theta_t = \theta_N + \Delta \left(\frac{e^{(1-\frac{t-1}{N-1})^p} - 1}{e - 1} \right), \quad p \in \{0.125 \text{ (Boulware)}, 0.725 \text{ (linear)}, 4.0 \text{ (conceder)}\}.$$

Decision rules

Target selection

- Compute θ_t using the chosen schedule; define $\theta_t^{\text{eff}} = \max(\theta_t - \epsilon, B)$.
- Player1 "leads" the schedule: in round t it uses $\theta_{t+1}^{\text{eff}}$ (clamped at $t = N$); Player2 uses θ_t^{eff} .
- Targets remain in raw utility while constructing counter-offers; discounts are applied only when evaluating incoming offers against the target and BATNA.

Accept / Walk

- **Accept** iff the discounted value of the opponent's offer meets both the discounted target and the discounted BATNA.
- **Walk** if every legal counter-offer that clears the target would concede zero items to the opponent (such offers are filtered out). Practically, this handles the high-BATNA edge cases where there is no room to bargain/concede.

Counter-offer construction

1. **Mirror the opponent's keep.** If the last opponent offer was x , their keep vector is $k = q - x$ for total quantities q . We search for k_y solving

$$\min_{k_y \in \mathcal{K}} \sum_{j=1}^T |k_{y,j} - k_j| \quad \text{s.t. } u_i(k_y) \geq \theta_t^{\text{eff}}, 0 \leq k_{y,j} \leq q_j,$$

where \mathcal{K} is the set of keep vectors induced by legal actions.

2. A depth-first search enumerates counts that satisfy the utility floor while minimising the ℓ_1 distance, ties go to the candidate closest to the target, and then to higher utility.
3. We're using the opponent's last offer as an opponent model, where we are trying to make them an offer that is as close as possible to their prior offer (in terms of the items they kept) while also adhering to our target.

We instantiate the Aspiration agent using the polynomial family with a maximum aspiration ceiling of 85% of the surplus above the BATNA. Acceptance slack is set to 5%, meaning the agent will accept offers slightly below its target threshold. The initial aspiration scale and concession exponent are coupled to the discount factor: patient games ($\gamma = 0.98$) use a Boulware profile that holds out early and concedes slowly, while impatient games ($\gamma = 0.9$) use a concedes profile that makes aggressive early concessions.

M BGL0

| Agent | BGL0 (LLM-only) | | | | |
|------------------|----------------------|-------------------|-------------------|-------------------|-------------------|
| | MENE Regret | NW% | NW ⁺ % | UW % | EF1 % |
| sonnet-3.7-c-5 | 33.24±0.77 | 87.95±0.05 | 53.33±0.10 | 85.87±0.05 | 36.35±0.13 |
| sonnet-3.7-c-6 | 23.37±0.49 | 89.61±0.09 | 55.94±0.11 | 86.99±0.06 | 36.29±0.18 |
| sonnet-3.7-r-c-0 | 20.39±0.61 | 89.80±0.05 | 50.49±0.11 | 85.58±0.04 | 43.86±0.14 |
| gem-2.0-f-c-2 | 33.61±0.73 | 85.94±0.10 | 45.34±0.13 | 82.65±0.07 | 32.68±0.16 |
| gem-2.0-f-c-5 | 42.62±0.69 | 85.66±0.10 | 47.61±0.14 | 84.42±0.07 | 35.80±0.35 |
| 4o-c-4 | 20.65±0.51 | 90.77±0.06 | 52.77±0.13 | 85.53±0.05 | 41.59±0.19 |
| 4o-c-5 | 7.33±0.36 | 91.99±0.07 | 54.93±0.14 | 87.08±0.05 | 39.62±0.20 |
| 4o-c-6 | 10.46±0.49 | 90.43±0.04 | 56.47±0.11 | 85.81±0.03 | 40.70±0.20 |
| o3-mini-c-0 | 2.12±0.21 | 92.37±0.06 | 63.08±0.18 | 87.16±0.05 | 55.96±0.21 |
| Agent | BGL0 (LLM+baselines) | | | | |
| | MENE Regret | NW% | NW ⁺ % | UW % | EF1 % |
| sonnet-3.7-c-5 | 42.84±1.32 | 86.48±.18 | 53.11±.27 | 85.08±.07 | 35.17±.22 |
| sonnet-3.7-c-6 | 31.79±1.35 | 88.56±.21 | 55.51±.29 | 86.58±.09 | 35.76±.24 |
| sonnet-3.7-r-c-0 | 23.28±.65 | 88.75±.16 | 50.17±.26 | 85.00±.07 | 43.03±.24 |
| gem-2.0-f-c-2 | 36.14±.74 | 85.23±.11 | 44.70±.25 | 82.17±.07 | 32.82±.20 |
| gem-2.0-f-c-5 | 48.14±.92 | 84.68±.16 | 46.29±.27 | 83.80±.09 | 35.13±.37 |
| 4o-c-4 | 24.08±.56 | 90.22±.11 | 51.82±.28 | 85.21±.07 | 40.17±.28 |
| 4o-c-5 | 9.16±.50 | 91.26±.12 | 53.94±.29 | 86.49±.09 | 38.42±.27 |
| 4o-c-6 | 12.25±.56 | 90.20±.10 | 54.85±.29 | 85.75±.07 | 39.13±.26 |
| o3-mini-c-0 | 1.70±.20 | 92.14±.11 | 61.89±.36 | 86.98±.08 | 54.27±.32 |
| soft | 51.58±1.36 | 85.41±.27 | 36.94±.34 | 82.83±.16 | 24.72±.34 |
| tough | 117.33±1.69 | 68.36±.17 | 2.66±.04 | 72.65±.10 | 0.03±.00 |
| walk | 93.54±.67 | 72.01±0.00 | 0.00±0.00 | 72.09±0.00 | 0.00±.00 |

Table 15: Bootstrap means (\pm standard errors) of regret, normalized Nash and utilitarian welfares, and EF1 frequencies, the last 3 expressed as percent points, for our LLM-only and LLM+baselines analyses on BGL0.