

Choose Your Agent: Tradeoffs in Adopting AI Advisors, Coaches, and Delegates in Multi-Party Negotiation

Kehang Zhu
Harvard University
Cambridge, United States

Nithum Thain
Google DeepMind
Toronto, Canada

Vivian Tsai
Google DeepMind
Mountain View, United States

James Wexler
Google DeepMind
Cambridge, United States

Crystal Qian
Google DeepMind
New York, United States

ABSTRACT

As AI usage becomes more prevalent in social contexts, understanding agent-user interaction is critical to designing systems that improve both individual and group outcomes. We present an online behavioral experiment ($N = 243$) in which participants play three multi-turn bargaining games in groups of three. Each game, presented in randomized order, grants *access to* a single LLM assistance modality: proactive recommendations from an *Advisor*, reactive feedback from a *Coach*, or autonomous execution by a *Delegate*; all modalities are powered by an underlying LLM that achieves superhuman performance in an all-agent environment. On each turn, participants privately decide whether to act manually or use the AI modality available in that game. Despite preferring the *Advisor* modality, participants achieve the highest mean individual gains with the *Delegate*, demonstrating a preference-performance misalignment. Moreover, we find suggestive evidence that delegation generates positive externalities: even non-adopting users in *access-to-delegate* treatment groups trend toward higher individual surplus, consistent with receiving higher-quality offers. Mechanism analysis reveals that the *Delegate* agent acts as a market maker, injecting rational, Pareto-improving proposals that restructure the trading environment. Our research reveals a gap between agent capabilities and realized group welfare. While autonomous agents can exhibit super-human strategic performance, their impact on realized welfare gains can be constrained by interfaces, user perceptions, and adoption barriers. Assistance modalities should be designed as mechanisms with endogenous participation; adoption-compatible interaction rules are a prerequisite to improving human welfare with automated assistance.

KEYWORDS

human-AI interaction, negotiation, delegation, LLM agents, behavioral experiment, multi-party

ACM Reference Format:

Kehang Zhu, Nithum Thain, Vivian Tsai, James Wexler, and Crystal Qian. 2026. Choose Your Agent: Tradeoffs in Adopting AI Advisors, Coaches, and Delegates in Multi-Party Negotiation. In *Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026)*, Paphos, Cyprus, May 25 – 29, 2026, IFAAMAS, 18 pages.

Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026), C. Amato, L. Dennis, V. Mascardi, J. Thangarajah (eds.), May 25 – 29, 2026, Paphos, Cyprus. © 2026 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). This work is licensed under the Creative Commons Attribution 4.0 International (CC-BY 4.0) licence.

1 INTRODUCTION

Large language models (LLMs) are shifting from passive tools to autonomous *agents*, capable of navigating complex social tasks and reshaping the incentives of multi-party environments. Understanding the design and effects of such systems is crucial for the emerging “agentic economy” [56, 60, 65]. A central challenge in this new economy is the tension between *control* and *delegation*: specifically, the normative question of when users *should* offload decision-making, and the behavioral question of when they actually *do*.

This tension is increasingly operationalized through a diverse spectrum of assistance modalities. Beyond traditional chat interfaces that passively respond to user input, emerging systems range from proactive assistants that initiate guidance [27, 46] to autonomous agents empowered with independent execution [45]. However, the impact of these design choices on strategic outcomes—and their scalability in multi-party environments—remains significantly underexplored.

Prior research on human-AI collaboration has primarily evaluated single-user, single-AI interactions in domains such as medical diagnosis, credit assessment, and risk prediction [2, 8, 28], often targeting objectives with verifiable ground truths [67]. These studies have established foundational insights, such as a preference-performance misalignment, where users may prefer sub-optimal agents that require less cognitive load [8, 14, 15].

Such human-AI interaction patterns become increasingly complex in social, real-world applications, which are defined by strategic interdependence, dynamic equilibria, and collective externalities. Crucially, while preference for higher-control AI modes has been observed in coding and classification tasks, its consequences in *multi-party* settings are qualitatively different: one user’s choice of modality affects the payoffs of others who never interact with AI at all. Recent work on autonomous LLM negotiations demonstrate that agent-to-agent bargaining can exhibit distinct risks and behaviors relative to human negotiation [52, 71]. A complementary line of research studies *principal-agent* interaction patterns in social contexts, where humans customize an agent (often via prompt writing) to negotiate autonomously on their behalf [32, 68]. While these studies establish the feasibility of autonomous negotiation,

they abstract away from the practical design choice of user controllability — specifically, how different modalities of interaction influence adoption and equilibrium outcomes.¹

We study how different allocations of agency affect behavior and welfare in a strategically interdependent setting, operationalizing the spectrum of AI agency into three distinct interaction modalities while holding the underlying model capability constant: an *Advisor* (proactive recommendations), a *Coach* (reactive feedback), and a *Delegate* (autonomous action).

We evaluate these modalities through a bargaining game [52] designed to evaluate AI capabilities in group negotiations, a setting with empirically exhibited human inefficiencies in trading behaviors [11, 57]. Our LLM agents utilize prompt scaffolding based on Gemini-2.5 Flash [17] that empirically achieves superhuman performance, allowing us to isolate the effects of the interaction structure. We then conduct a randomized, within-subjects experiment ($N = 243$) where participants engage in three successive games, using each modality in a counterbalanced order.

Our study makes the following contributions:

- **A micro-economy design for testing agentic assistance.** We introduce and run a randomized, within-subject experiment in a strategically interdependent three-player bargaining game that cleanly varies *modalities* of LLM integration (Advisor, Coach, Delegate) while holding underlying model capability fixed.
- **Evidence of preference vs. welfare under endogenous use.** We document preference–performance misalignment in a strategic, interdependent setting: participants prefer the higher-control *Advisor* interface, yet groups achieve the highest average surplus when *Delegate* access is available.
- **Externalities and a market-making mechanism.** We find evidence of spillover effects from *Delegate* usage, driven by a concrete mechanism: Delegate agents propose larger, asymmetric, Pareto-improving trades that shift the distribution of available deals for all players—a market-making role that benefits non-users even when they never invoke the AI. This mechanism is confirmed via distributional tests on accepted trade surplus (Kolmogorov–Smirnov, $p = .002$).

Although the LLM agents exhibit measurable superhuman capabilities in this setting, our randomized study shows that users do not reliably take them up due to interaction frictions. This suggests that providing super-human agentic capabilities alone can be insufficient for improving human outcomes; therefore, interface design and interaction patterns are central to realizing benefits in collective systems. By documenting this gap and presenting empirical evidence on user preferences, individual and aggregate welfare, and spillover effects, our work offers an early baseline for studying human–agent interaction in strategic, multi-party ecosystems.

2 RELATED WORK

2.1 Negotiation and Group Decision Making

LLM-based agents demonstrate increasingly strong social capabilities across negotiation, persuasion, mediation, and multi-agent

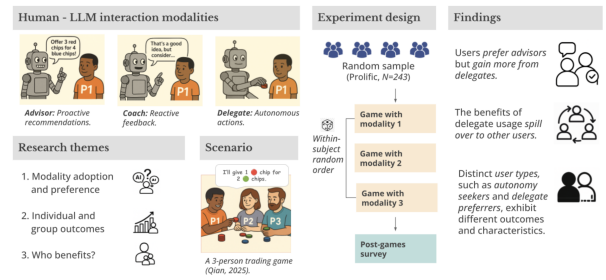


Figure 1: Overview of the experimental design and contributions. Participants ($N=243$) engaged in three-person bargaining games with access to three LLM assistance modalities: *Advisor* (proactive recommendations), *Coach* (reactive feedback), or *Delegate* (autonomous actions) — with within-subject game randomization. We find evidence of a preference-performance misalignment, positive spillover effects in the Delegate games, and heterogeneous outcomes for different participant types.

coordination tasks [1, 5, 31, 40, 41, 59, 61, 62, 64]. Modern LLMs have exhibited near-human-level strategic performance in complex multi-party settings such as *Diplomacy* [24], outperform human debaters in randomized tournaments [61], and even produce moral judgments preferred by human evaluators [48].

Beyond individual capabilities, LLMs have been deployed to support collective processes, such as facilitation, moderation, and mediation, in committees, communities, and group deliberation [7, 20, 23, 34, 64]. However, most human-AI studies evaluate individual users (i.e., 1:1 human-AI dyads) in tasks with ground-truth data, such as diagnosis, lending, or deception detection [2, 3, 28, 30, 36, 37]. Moving from dyads to multi-party bargaining introduces additional complexities as one participant’s actions change the opportunity set faced by others, creating structural externalities.

To enable controlled comparison in such environments, Qian et al. [52] propose a stylized multi-party bargaining game with induced values, abstracting away subjective goals [18] while preserving strategic interdependence. Related work also investigates fully autonomous LLM negotiators through prompt-designed competitions [32, 68] and agent–agent market simulations [71]. While these studies primarily evaluate agent capabilities in bargaining environments, our study evaluates the human takeup of such capabilities through varying interaction modalities.

2.2 Delegation in Human–AI Interaction

Full delegation to agentic systems remains relatively novel in empirical studies of human–AI teaming; a recent meta-analysis [67] reviewed over 100 experiments and found that only a small subset involve structured delegation of decision authority. When delegation is studied, it typically appears in isolated tasks where complementarity can be engineered; for example, hybrid human–AI assessments in physical therapy [39], post-editing pipelines in summarization [38], or selective triage allocations based on uncertainty or rule-based heuristics [4, 54].

¹For example, coding agents like Claude-Code offer three distinct modes of assistance: *planning mode*, *ask before edits*, and *edit automatically*.

In group settings, agent delegation can increase cooperation in collective-risk dilemmas [23]. AI-mediated group consensus can be perceived as clearer and less biased [64]. In multi-party negotiation, AI moderation can raise fairness and efficiency perceptions [34].

2.3 Adoption and Control in AI Assistance

Despite potential benefits of AI adoption and delegation, take-up of such systems by human users often involves a calculated tradeoff between outcome quality and perceived control [47, 58]. Human-AI teams may fail to achieve synergy even when AI improves accuracy [28, 67], in part because designs that enhance objective performance may increase cognitive effort or monitoring costs [14, 15, 19, 51, 63]. Users must also maintain accurate mental models of when AI is likely to err [8]; miscalibrated confidence or overestimation of personal ability [35, 69] can lead to rejection of useful assistance. *Algorithm aversion* further reduces adoption after observing small errors [21, 22], though allowing human-in-the-loop intervention can restore perceived agency.

In sum, these studies emphasize that adoption should not be treated as a proxy for effectiveness: users choose workflows that balance minimizing cognitive load and maximizing objective performance [25, 49]. Our work connects these insights to a strategic, multi-user environment: we examine how different human-AI interaction modalities shape adoption behavior, individual outcomes, and spillovers that arise when participants’ decisions are mutually interdependent.

3 EXPERIMENT

We frame our empirical inquiry through three sets of research questions. Broadly, when access to AI is given to individuals in a multi-user context, who uses it, what do they gain, and who benefits the most?

RQ1: Welfare effects. At the group level, does access and adoption of LLM assistance—in any of the three modalities (*Advisor*, *Coach*, *Delegate*)—increase *individual*- and *group*- level surplus, relative to a no-AI baseline?

RQ2: Spillovers and externalities. Do any of the assistance modalities create positive or negative externalities for non-users within the same group?

RQ3: Adoption and preferences. How do participants’ stated preferences and usage decisions vary across the *Advisor*, *Coach*, and *Delegate* modalities, independent of their welfare effects?

3.1 Game Setting

We explore the negotiation dynamics in the context of a chip bargaining game introduced by Qian et al. [53], where participants take turns exchanging chips of different colors, with randomly assigned private valuations, with the goal of maximizing their surplus. During each turn, one participant *proposes* an offer (e.g., 7 red for 3 green chips), and the other two simultaneously and privately decide whether to *accept* or *decline*. If both accept, one is chosen randomly to clear the trade. The game ends after nine turns; participants leave with any surplus earned beyond the initial value of their chips.

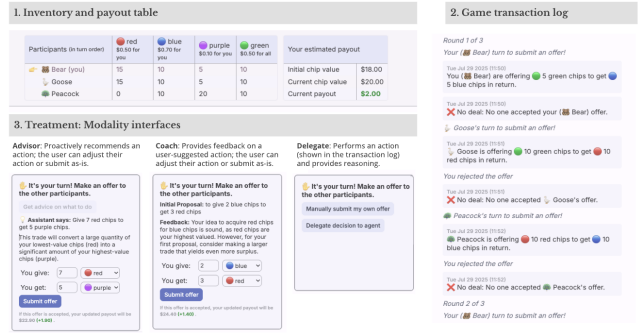


Figure 2: Relevant game components. Panels 1. and 2. show interface properties visible across all modalities. Panel 3 illustrates the proposal generation interface shown in each of the treatment modalities. More detailed figures of game interfaces are provided in Appendix E.

This environment provides several properties that enable our investigation. First, by endowing participants with pre-defined preferences, it provides an objective ground truth for performance and outcome measurement, which is a requirement in empirical trading games exploring similar behaviors [12, 33, 55]. Second, the combination of information asymmetry and restricted communication create a setting that necessitates strategic reasoning.² Finally, the multi-player, non-zero-sum design allows for the direct measurement of group-level externalities, enabling the measurement of RQ1 (individual and group outcomes).

Game-theoretical benchmark. Furthermore, the structure of the game allows computation of a Pareto-efficient benchmark for total surplus. Qian et al. [53] measures the *performance*, or surplus gain, as the surplus achieved by a group or individual divided by the maximum possible (Pareto-efficient) surplus for that specific game configuration, allowing for normalized comparison across different games and conditions. Throughout the paper, we refer to these measures as *scaled group surplus* and *scaled individual surplus*.

All-Human baseline. For our baseline, we reference the human-only performance established in Qian et al. [53]. To ensure comparability, we verified that the recruitment criteria, payout rates (\$10 base + bonus), and interface mechanics were identical to the current study. They found that human-only groups ($N = 72$) achieve a mean scaled group surplus of $0.537 (\pm 0.024)$. This suboptimal performance was attributed to human tendencies toward conservative trades and a “fairness” norm (e.g., 1-for-1 swaps), which systematically limited collective gains.

3.2 Conditions

We compare three treatment arms — *Access-to-Advisor*, *Access-to-Coach*, and *Access-to-Delegate*— against two no-AI baselines from Qian et al. [52]: an all-human baseline and a *theoretical optimum* baseline (a linear-programming solution under full public valuations) that benchmarks optimized performance.

²The complexity of this game prevents any dominant strategies that can be calculated ex-ante [53].

The intervention is *access* to a given treatment modality; in each game, each player can choose between taking action independently of AI assistance, or using the available AI intervention. All players within the game have access to the same intervention modality (e.g., a *Delegate* game). A player can only interact with one modality per game and can choose to use assistance across each turn in the game; on each of the nine turns, they independently decide whether to invoke the available AI modality or act manually:

- **Advisor:** The advisor agent proactively recommends an action, either an offer or an accept/reject response, and provides a rationale. The user can accept or revise this recommendation before submitting it.
- **Coach:** The user first composes their intended action. The coach agent provides feedback on the user’s plan, which the user can incorporate or disregard.
- **Delegate:** The delegate agent autonomously generates and executes an action on the user’s behalf. The user cannot veto or modify the AI’s decision, but can view the agent’s rationale.

Figure 1 visualizes the experiment design. First, participants are introduced to the game rules and AI modalities, completing mandatory comprehension checks to ensure understanding. The instructions intentionally avoided framing the AI modalities as “superhuman” to prevent biasing towards delegation. The participants completed a pre-game survey and played three games in succession, each with access to a different AI assistance modality (Advisor, Coach, Delegate) in randomized order. Finally, the participant completed a post-game survey.

3.3 LLM Agent Design and Implementation

All assistance modalities were powered by an LLM-based agent, developed with two design goals in mind: i) its capabilities in the bargaining game must exceed human baselines, and ii) it needs to respond to the user with minimal latency. Our agent used the commercially-available Gemini-2.5-Flash API, with an output token limit of 8,192, thinking token budget of 2,048, and a throughput limit of 1,000 tokens/second to ensure responsive assistance. For robustness, model queries were implemented with two retries and a fallback mechanism that reverted users to manual override mode when necessary. Empirically, this implementation yielded a 99.71% assistance success rate across 2,519 queries, with manual fallback triggered only seven times.

While more powerful models like Gemini 2.5-Pro achieved slightly higher performance in simulations, the >40-second response latency was impractical for feedback during a real-time game. Gemini-2.5-Flash provided an average response time of 10.55 seconds, meeting the usability threshold for maintaining user engagement [42, 44]. All three assistance modes were implemented as lightweight scaffolds on top of the same underlying agent, with supplemental reasoning text.³

All-AI baseline. We developed prompt scaffolding that outperformed the baseline agents in Qian et al. [53], which exhibited human-level capabilities. In simulation, our agent achieved a scaled

surplus of 0.595 (± 0.024), significantly outperforming the human baseline (0.537 ± 0.024).

3.4 Design and Analysis

We utilize a within-participant design. Each participant played three separate games with access to a mode in randomized order, alongside two other players assigned to the same mode. Because modality order is counterbalanced across participants, any learning effects from successive games are distributed equally across conditions and do not confound modality comparisons.

We collected the following performance measures:

- **Individual/group surplus gain:** participant’s or group’s surplus change relative to the original chip values.
- **Proposal and decisions:** trading proposals and other players’ responses (reject/accept).
- **AI takeover per turn:** whether the participant used AI assistance at certain turn.

We also collected self-reported subjective measures.

Pre-game survey. To inform RQ3, we employed a pre-game survey to capture user attributes.⁴ This survey captured three constructs on a 5-point Likert scale:

- **Trust in AI:** Perceived ability, trustworthiness, insight, and helpfulness [26]. For analysis, we compute a composite pre-game trust score as the simple average of these four items.
- **Prior expertise:** Prior familiarity in games or tasks similar to the bargaining game; the majority of participants reported low prior familiarity, suggesting results reflect novice-to-intermediate users rather than domain experts.
- **Confidence:** Confidence in their ability to play the bargaining game well.

Post-game survey. Participants completed a survey after each of the three game rounds, and a final comparative survey. These surveys aimed to capture the following constructs:

- **Satisfaction:** Satisfaction with final trading outcomes.
- **Mental effort:** Cognitive load exhibited in the bargaining games.
- **Preference:** Choice of which AI modality they would prefer to use for future games.

Our data was analyzed using mixed-effects models. The choice of game modality was modeled as a fixed effect, and group as a random effect to account for the fact that three participants form a group (i.e., all the measurements were not statistically independent). We also report Holm-Bonferroni corrections to account for multiple comparisons [29] noting that this may lead to false negative results due to the conservative nature of the correction [16, 43].

3.5 Procedure and Participants

The game interface was implemented and deployed using Deliberate Lab [66], an open-source experimentation platform.⁵ 324 participants were recruited from the Prolific recruitment platform under an IRB-approved protocol, with no additional selection criteria [50].

³The full prompts for these agents are in Appendix F, including example generated proposals for each modality, and screenshots from the interfaces are in Figure 2.

⁴A complete list of pre- and post- game survey questions and responses is provided in Appendix D.

⁵Additional game interface implementation details are provided in Appendix E.

Our final sample includes $N = 243$ participants over 81 groups of three,⁶ with 13 to 15 groups in each of the six unique orderings across the three treatments, involving 6, 561 trading decisions. Participants received a \$10.00 base payment, plus a performance-based bonus that averaged \$4.50, for approximately 56.4 minutes of their time. The bonus, calculated as the average individual surplus across the three games, was designed to align participant incentives with surplus maximization.

4 RESULTS

4.1 LLM Delegation leads to the highest gains (RQ1)

To evaluate the impact of AI assistance on negotiation outcomes, we employed Linear Mixed-Effects Models (LMM) to account for the nested structure of repeated measures within negotiation groups. The general formulation of the model for group-level efficiency is defined as:

$$Y_{ij} = \beta_0 + \beta_1 \cdot \text{Condition}_{ij} + u_j + \epsilon_{ij} \quad (1)$$

Where:

- Y_{ij} represents the scaled surplus for session i in group j .
- β_0 is the fixed intercept (representing the Human Baseline).
- β_1 represents the vector of fixed effects for the AI modalities (Delegate, Advisor, Coach).
- $u_j \sim \mathcal{N}(0, \sigma_u^2)$ is the random intercept for group j , accounting for group competencies.
- $\epsilon_{ij} \sim \mathcal{N}(0, \sigma_\epsilon^2)$ is the residual error.

We applied Holm-Bonferroni corrections to the family of pairwise comparisons between the AI modalities and the Human Baseline to control for family-wise error rates.

Group-level surplus (RQ1a). **Delegation** was the only modality to show a positive directional trend in group welfare compared to the human benchmark, though this did not survive multiple-testing correction ($p_{adj} = .10$). Groups with access to the Delegate agent achieved significantly higher scaled surplus than the Human Baseline ($\beta = 0.084, SE = 0.040, p = .034$). However, due to the conservative nature of the Holm-Bonferroni correction applied across all three modalities, this difference was rendered marginally non-significant ($p_{adj} = .10$). In contrast, neither the **Advisor** ($\beta = 0.006, p_{adj} = 1.0$) nor the **Coach** ($\beta = 0.026, p_{adj} = 1.0$) conditions showed any significant improvement over the baseline, nor did they exhibit the positive directional trend observed in the Delegate condition.

Individual-level surplus (RQ1b). A similar pattern emerged at the individual level in Table 2. We updated the model to include individual random effects nested within groups. While the Delegate mode showed a positive coefficient for individual surplus ($\beta = 0.028$), it did not reach statistical significance after correction ($p_{adj} = .20$). The Advisor and Coach modes showed negligible differences from the baseline. These results suggest that while the Delegate mode introduces efficiency gains into the market, the variance in individual

⁶Because the design required full group participation across all three games, we excluded any group in which one or more members failed to complete the session. Excluded participants did not differ significantly from completers on pre-game trust or prior expertise ($p > .10$), suggesting no systematic selection bias.

Parameter	Coef.	Std. Error	z-value	P-value	P_{adj}
Intercept	0.537	0.032	16.737	0.000	—
Advisor (vs Human)	0.006	0.040	0.141	0.888	1.000
Coach (vs Human)	0.026	0.040	0.665	0.506	1.000
Delegate (vs Human)	0.084	0.040	2.127	0.033*	0.100
Group Variance	0.006	0.017	—	—	—

Table 1: Mixed Linear Model Results for RQ1a (Group-Level Efficiency)

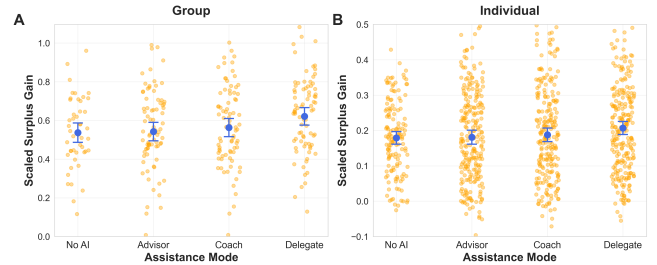


Figure 3: Access to delegation increases both group and individual surplus. Panel (A) shows the comparison of scaled group surplus gain. Panel (B) shows the comparison of scaled surplus gain for individuals.

outcomes—likely driven by heterogeneous adoption rates—dilutes the statistical signal at the individual level compared to the group level.

4.2 Delegation Yields Positive Spillover Effects (RQ2)

We compared the performance of non-users in groups where AI was available against the Human Baseline. We define a *non-user* as a participant who did not use the AI for proposals *in that specific condition*, whereas a user is defined as someone who utilized AI assistance at least once for a proposal. This definition echoes prior findings that the act of making a proposal drives the final surplus significantly more than merely answering an offer, which is less strategic [52]. We also tested a Linear Mixed-Effects Model treating each players' own AI usage in proposal making as a continuous variable.⁷

4.2.1 Model Specification. We modeled the individual scaled surplus using a Linear Mixed-Effects Model (LMM) to account for group-level dependencies. The data subset included all participants from the Human Baseline and identified non-users from the AI conditions. The model is defined as:

⁷Results are detailed in Appendix A.

Parameter	Coef.	Std. Error	z-value	P-value
Intercept	0.179	0.012	14.757	0.000
Advisor (vs Human)	0.002	0.015	0.122	0.903
Coach (vs Human)	0.009	0.015	0.573	0.567
Delegate (vs Human)	0.028	0.015	1.832	0.067

Table 2: Mixed Linear Model Results for RQ1b (Individual-Level Efficacy)

$$\begin{aligned}
Y_{ij} = & \beta_0 + \beta_1 \cdot \mathbb{I}(\text{Condition}_{ij} = \text{Delegate_NonUser}) \\
& + \beta_2 \cdot \mathbb{I}(\text{Condition}_{ij} = \text{Advisor_NonUser}) \\
& + \beta_3 \cdot \mathbb{I}(\text{Condition}_{ij} = \text{Coach_NonUser}) + u_j + \epsilon_{ij}
\end{aligned} \quad (2)$$

Where:

- β_0 still represents the intercept (mean surplus of the Human Baseline).
- $\beta_1, \beta_2, \beta_3$ represent the marginal effect of being a non-user in the Delegate, Advisor, and Coach conditions, respectively.

Parameters	Coef.	Std. Err.	z-value	P-value	Preferred Mode	Number of Participants (N)	Share (%)
Intercept	0.179	0.012	14.635	0.000	Advisor	107	44.0
Delegate Non-User (vs Human)	0.039	0.020	1.926	0.054	Coach	37	15.2
Advisor Non-User (vs Human)	-0.017	0.022	-0.768	0.443	Delegate	47	19.3
Coach Non-User (vs Human)	0.012	0.020	0.627	0.530	None	52	21.4

Table 3: Linear Mixed-Effects Model results comparing the individual surplus of non-users in AI conditions against the Human Baseline.

4.2.2 Results. Table 3 presents results from the regression analysis. We observed a notable positive trend for non-users in the **Delegate** condition. Delegate non-users achieved a mean individual scaled surplus of 0.218, which is 21.6% higher than the Human Baseline ($Mean = 0.179$). The LMM analysis showed a positive coefficient for Delegate non-users ($\beta = 0.039, SE = 0.020$), which approached statistical significance in the raw comparison ($p = 0.054$). However, after applying the Holm-Bonferroni correction for the three comparisons, this result did not meet the standard threshold for significance ($p_{adj} = 0.16$). Non-users in the **Advisor** ($Mean = 0.162$) and **Coach** ($Mean = 0.191$) conditions did not show significant improvements over the baseline (Advisor: $\beta = -0.017, p = 0.44$; Coach: $\beta = 0.012, p = 0.53$).

While the evidence for RQ2 does not reach statistical significance after multiple testing corrections, the descriptive data reveals a compelling trend: the Delegate modality was the unique condition in which non-users achieved a higher mean surplus than both the human baseline and the active AI adopters within the same condition (0.218 vs. 0.201).

4.3 Users Prefer the Advisor Modality (RQ3)

Participants strongly preferred the Advisor modality—a finding that directly contradicts their objective performance ($\chi^2 = 45.03, p <$

.001). We analyzed participants’ responses to the multiple-choice survey question: “If you were to play again, which AI assistance mode would you prefer to use?”. Previous literature suggests that users prefer higher-control modalities; our results confirm this pattern while extending it to a multi-party strategic setting where the preference carries collective costs.

4.3.1 Results. A Chi-Square test revealed significant differences in user preferences across the three AI modalities ($\chi^2 = 45.03, p < .001$). To identify specific drivers of this preference, we conducted pairwise comparisons with Holm-Bonferroni corrections. The results provide **partial support** for RQ3, revealing a strong preference for the Advisor but not for the Coach as shown in Table 4:

- **Advisor vs. Delegate:** Consistent with RQ3, users significantly preferred the Advisor mode ($N = 107$) over the Delegate mode ($N = 47$) ($p_{adj} < .001$).
- **Advisor vs. Coach:** Users also overwhelmingly preferred the Advisor over the Coach ($N = 37$) ($p_{adj} < .001$).
- **Delegate vs. Coach:** There was no significant difference in preference between the Delegate and the Coach ($p_{adj} = .326$).

Table 4: Preference differences across agent types: users dominantly prefer the advisor modality.

4.3.2 Rationales for preferences. We furthermore conducted a semantic analysis of open-ended rationale text [13], identifying four recurring themes for preference: (i) *trust and control*, (ii) *ease of use and cognitive offloading*, (iii) *effectiveness and performance*, and (iv) *other*.⁸

Advisor-preferrers frequently credited the AI’s effectiveness and performance as the reason for their selection; many also referred to trust and control and ease of use. From P92:

“I love the advisor. It helps when you get into the weeds of the game when the strategies become less obvious. I also like that I still have full control.”

Participants who preferred none of the modalities (“**autonomy-seekers**”) also cited trust and control, but as grounds for disengagement. From P96:

“I don’t trust AI bots; I feel I can make better decisions on my own.”

Delegate-preferrers valued ease of use and cognitive offloading. From P47 and P99:

“Delegate helped with ease of decision making and made it easiest for me.”

“I prefer that someone else make the decision.”

⁸Thematic analysis methodology provided in Appendix D.1.

Coach-preferrers highlighted effectiveness and performance as the main rationale for their selection. From P266:

“Coach helped me see things that I didn’t see myself like a real coach.”

4.4 Explorative Analysis

4.4.1 Participants prefer AI assistance for generating offers over responding. Table 8 shows within-participant AI usage rates; across all modalities, participants used AI assistance significantly more for offer generation than for offer response.

4.4.2 AI usage over time remains stable for offers, but declines for responses. As shown in Figure 5, participants’ reliance on AI for offer generation remained relatively steady across the three rounds of play, whereas usage for responses declined significantly in all modalities (coef. = -0.407, $p < 0.001$).

4.4.3 Acceptance of AI suggestions differs sharply by modality. Requesting AI assistance did not guarantee adoption of its suggestions. The Delegate modality yielded a 100% acceptance rate by design. However, in Advisor mode, users sent the Advisor’s recommended offer without modification only 70.6% of the time. In Coach mode, users retained their initial offer 69.5% of the time and their initial response 96% of the time, even when the AI recommended a different or opposing action (e.g. rejecting rather than accepting an offer). Higher **pre-game trust** in AI predicts greater AI usage across all modes ($r \sim .25$, $p < .01$).

4.4.4 Post-game Mental load and Pre-game Trust. Post-game mental load predicts lower performance and greater AI reliance. An OLS regression shows that higher self-reported mental effort (i.e., perceiving the game as more difficult) was associated with lower average surplus across all games (coef = -0.012, $p = 0.023$). Participants who preferred any AI mode reported 20% higher mental effort than autonomy-seekers ($p < .01$; Table 9).

5 DISCUSSION

5.1 Decomposing Delegate Advantages

The *access-to-delegate* yielded an increase in scaled surplus of 2.8 percentage points relative to the Human Baseline (No AI). To understand the drivers of this performance lift, we conducted a contribution decomposition analysis separating direct user benefits from environmental externalities.

$$\Delta_{\text{Total}} = \underbrace{p \cdot (\bar{Y}_{\text{User}} - \bar{Y}_{\text{Base}})}_{\text{Internal Effect (Direct)}} + \underbrace{(1 - p) \cdot (\bar{Y}_{\text{Non}} - \bar{Y}_{\text{Base}})}_{\text{External Effect}}$$

5.1.1 Decomposition of welfare gains. The total welfare advantage of the Delegate mode is derived from two distinct sources:

- (1) **Direct contribution (internal effect):** Approximately 51.2% of the total benefit is attributed to the improved outcomes of users who actively adopted the Delegate agent. However, statistical analysis reveals that the average improvement for these adopters (+0.022) was not statistically significant ($p = .148$).

- (2) **External contribution:** The remaining 48.8% of the gain stems from structural improvements in group dynamics that benefited non-adopters. Non-users in Delegate groups experienced a surplus gain of +0.039, where the exploratory unadjusted p-value is significant ($p = .033$).⁹

This decomposition shows that the benefits of delegation apply to both users and non-users. Delegation acts as a *market maker*, improving the trading environment for all participants.

5.1.2 The mechanism: high-value trade proposals. Next, we evaluate *why* non-users benefit significantly more when others delegate their proposal generation. The spillover effect arises not from higher rates delegate AI adoption, as acceptance rates do not significantly differ across modes (see Table 6), but from the quality of the offers generated by the Delegate agents.

The Delegate agent proposes larger, asymmetric trades that unlock higher-surplus, mutually-beneficial trades. This dynamic is visualized in Figure 4. Among all accepted offers, only those originating from a Delegate agent systematically shifted the distribution of the receiver’s surplus upward (Kolmogorov–Smirnov, $p = .002$). This indicates that Delegate agents propose more high-volume, high-gain agreements that benefit their human counterparts. While the proposer’s average gain in these trades was slightly lower—representing a strategic concession to ensure acceptance—the overall distribution of the proposer’s surplus also shifted significantly (Kolmogorov–Smirnov, $p = .030$), reflecting deals that created more aggregate value.

In contrast, in the Advisor and Coach modes, human intervention—such as editing suggestions or ignoring feedback—filtered out these high-value but complex offers, diluting the final trade distribution closer to those from the human-only baseline.

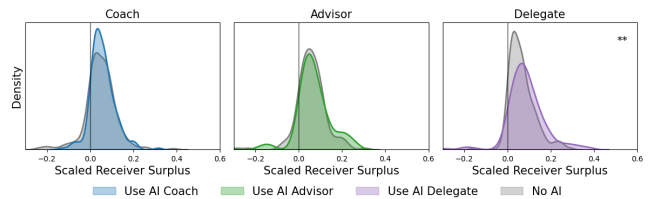


Figure 4: Receiver surplus distributions of accepted trades. Each panel compares the surplus of AI-assisted sender offers with non-AI-assisted offers. A significant upward shift in receiver surplus is observed only in the Delegate condition ($p < 0.01$), indicating that Delegate agents propose trades that are mutually beneficial.

5.2 Tradeoff Between Control, Cognitive Load, and Performance

Our results reveal a distinct *Preference-Performance Misalignment*. In RQ3, we find that participants strongly prefer the Advisor mode (43.7%) over the Delegate mode (19.1%). However, in RQ1, the Advisor provides no statistically significant welfare gain over the baseline, while the Delegate yields the highest objective economic value.

⁹We report both unadjusted and multiplicity-adjusted p-values $p_{\text{adjusted}} = 0.1$; we treat the decomposition evidence as exploratory.

A qualitative analysis of rationales suggest that *advisor-preferrers* are willing to sacrifice economic efficiency to maintain a sense of low-effort agency.

Cognitive load and assistance takeup. The “cognitive miser” hypothesis [15, 47], posits that users default to the path of least cognitive resistance. In our setup, the *Delegate* mode offered the ultimate cognitive offloading, automating the complex reasoning of proposing trades. Yet, users significantly favored the *Advisor*, which required them to actively evaluate and approve recommendations.

In our study, the tension appears to be between *agency* and *preference*. Users accepted the higher cognitive load of the *Advisor*, stating a preference for controllability (see participant rationales in Sec 4.3.2).

Coach games exhibited non-significant directionally higher group surplus (average 2.6 pp better than human baseline and 2 pp better the advisor group), though **RQ3** showed that users predominantly preferred *Advisor* assistance over *Coach*. This inverse relationship between user preference and objective utility aligns with the theory of *Cognitive Forcing Functions* Bućinca et al. [15]. The *Coach* mode requires users to first formulate a proposal before receiving feedback, inducing System II analytical thinking [?]. Conversely, the *Advisor* mode allows for passive reception of optimal moves, which may lead to superficial adoption without strategic internalization.

Capability calibration as an alternative mechanism. An alternative explanation for limited delegation is not a resistance to cognitive offloading, but uncertainty about the AI’s reliability and failure modes in a strategic, high-stakes setting. Our instructions intentionally avoided sharing with human participants that the agents exhibit superhuman AI performance (Section 3.1), to ensure that participants enter the game with uncalibrated beliefs about the agent’s capabilities. In such cases, reluctance to delegate can reflect capability uncertainty and algorithm aversion, rather than purely an intrinsic control premium.

As an exploratory check, pre-game trust predicts higher takeup of assistance (Sec 4.4.4). This interpretation aligns with our qualitative rationales: autonomy-seekers most often cite *Control & Trust* concerns, whereas delegate-preferrers emphasize *Ease of Use & Cognitive Offloading* (Appendix D.1). These findings validate that capability calibration is an important consideration for system design: interfaces that surface model confidence, show counterfactual outcomes, or provide a veto window for delegated actions may narrow the gap between objective gains and realized takeup [9, 67].

6 LIMITATIONS

6.1 Task Design and Ecological Realism

We intentionally use a stylized three-player chip-trading game as a contained environment for studying agentic assistance in *strategic, interdependent* decision making. The environment captures several properties that are central to real-world bargaining and resource allocation — private information, mixed-motive incentives, multi-party externalities, and hard budget constraints — while allowing precise control over payoffs and counterfactuals. This level of abstraction is consistent with work that uses games and simulations as testbeds for multi-agent behavior and mechanism design [6, 10, 70]. At the same time, it omits important real-world factors such as

unconstrained natural language communication, long-term relationships and reputation, and domain-specific norms [18] (e.g., labor negotiations, procurement, or humanitarian contexts). Our findings should therefore be interpreted as evidence about *patterns* of human-AI interaction under strategic interdependence, not as direct predictions for any specific deployed system.

6.2 Superhuman Model Capabilities

Our setup utilized LLM scaffolding that achieved superhuman performance. Not all consumer AI agents operate at this level, or operate in conditions with a measurable theoretical optimum.

6.3 Temporal Dynamics and Learning Effects

Our experiment measures outcomes from a single game per modality, with at most nine opportunities to interact with a given assistant (three proposals and six responses). This design provides clean comparisons across conditions, but is likely too short to capture how trust, reliance, and strategy evolve over time. Prior work on automation shows that trust is dynamic and history-dependent, shaped by sequences of successes and failures [26]. In longer horizons, users might either grow more comfortable with delegation as they observe its benefits, or react strongly to occasional failures and revert to manual control.

6.4 Interface and Design Variations

Our interface deliberately held presentation details constant across modalities (e.g., similar rationales, fixed framing, no explicit confidence scores) in order to focus on the structure of assistance. Prior work shows that explanation style, framing, and transparency can substantially influence both subjective trust and objective performance [9, 67]. Follow-up work could explore richer interface variations, such as offering users choice over dynamic or mixed modalities, providing “veto windows” for delegated actions, or surfacing model confidence and counterfactual outcome.

7 CONCLUSION

As AI systems increasingly assist users in collaborative contexts, from workplace coordination to online markets, they can change the dynamics of collaboration. While much attention has been paid to increasing these assistive capabilities, the *form* of human-AI interaction remains underexplored.

To address this, we conducted a within-participants randomized experiment in a multi-person bargaining game (N=243), comparing three assistance modalities — *Advisor*, *Coach*, and *Delegate* — while keeping the underlying model capabilities constant. We found a preference-performance misalignment: participants preferred the autonomy-preserving *Advisor* the most, yet achieved the highest payoffs, both individually and collectively, when using the fully autonomous *Delegate*. This advantage did not stem from higher adoption rates in the *Delegate* modality, but rather from the lack of user intervention that would dilute the model’s high-quality proposals. Delegation improved outcomes for adopters and also created positive externalities for peers, lifting group performance through “market-making” offers that improved the opportunity set for everyone in the game.

As LLM capabilities advance, a shift toward autonomous delegation may be increasingly likely, yet our findings surface a persistent tension between individual agency and systemic efficiency. Interfaces are not merely a user experience layer—they are part of the mechanism itself. Agentic systems deployed in social settings should therefore be evaluated as market-level interventions, with explicit attention to externalities. While this work provides an empirical baseline, further multidimensional analysis of preferences, market design, and human interaction is needed to narrow the gap between what users prefer, what they adopt, and what ultimately improves individual and collective welfare.

REFERENCES

- [1] Sahar Abdelnabi, Amr Gomaa, Sarath Sivaprasad, Lea Schönherr, and Mario Fritz. 2023. LLM-Deliberation: Evaluating LLMs with Interactive Multi-Agent Negotiation Games. (2023).
- [2] Nikhil Agarwal, Alex Moehring, Pranav Rajpurkar, and Tobias Salz. 2023. *Combining human expertise with artificial intelligence: Experimental evidence from radiology*. Technical Report. National Bureau of Economic Research.
- [3] Nikhil Agarwal, Alex Moehring, and Alexander Wolitzky. 2025. *Designing Human-AI Collaboration: A Sufficient-Statistic Approach*. Technical Report. National Bureau of Economic Research.
- [4] Nikhil Agarwal, Alex Moehring, and Alexander Wolitzky. 2025. *Designing Human-AI Collaboration: A Sufficient-Statistic Approach*. NBER Working Paper 33949. National Bureau of Economic Research. <https://doi.org/10.3386/w33949>
- [5] Saaket Agashe, Yue Fan, Anthony Reyna, and Xin Eric Wang. 2023. Llm-coordination: evaluating and analyzing multi-agent coordination abilities in large language models. *arXiv preprint arXiv:2310.03903* (2023).
- [6] Mohammed Alsobay, David G Rand, Duncan J Watts, and Abdullah Almaatouq. 2025. Integrative Experiments Identify How Punishment Impacts Welfare in Public Goods Games. *arXiv preprint arXiv:2508.17151* (2025).
- [7] Mohammed Alsobay, David M Rothschild, Jake M Hofman, and Daniel G Goldstein. 2025. Bringing Everyone to the Table: An Experimental Study of LLM-Facilitated Group Decision Making. *arXiv preprint arXiv:2508.08242* (2025).
- [8] Gagan Bansal, Besmira Nushi, Ece Kamar, Walter S. Lasecki, Daniel S. Weld, and Eric Horvitz. 2019. Beyond Accuracy: The Role of Mental Models in Human-AI Team Performance. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing (HCOMP)*, Vol. 7. 2–11.
- [9] Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. 2021. Does the whole exceed its parts? the effect of ai explanations on complementary team performance. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–16.
- [10] Federico Bianchi, Patrick John Chia, Mert Yuksekgonul, Jacopo Tagliabue, Dan Jurafsky, and James Zou. 2024. How well can llms negotiate? negotiationarena platform and analysis. *arXiv preprint arXiv:2402.05863* (2024).
- [11] Ken Binmore, Ariel Rubinstein, and Asher Wolinsky. 1986. The Nash bargaining solution in economic modelling. *The RAND Journal of Economics* (1986), 176–188.
- [12] Olivier Bochet, Manshu Khanna, and Simon Siegenthaler. 2024. Beyond dividing the pie: Multi-issue bargaining in the laboratory. *Review of Economic Studies* 91, 1 (2024), 163–191.
- [13] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative Research in Psychology* 3, 2 (2006), 77–101. <https://doi.org/10.1191/1478088706qp063oa>
- [14] Zana Bućinca, Phoebe Lin, Krzysztof Z Gajos, and Elena L Glassman. 2020. Proxy tasks and subjective measures can be misleading in evaluating explainable AI systems. In *Proceedings of the 25th International Conference on Intelligent User Interfaces*. 454–464. <https://doi.org/10.1145/3377325.3377498>
- [15] Zana Bućinca, Maja Barbara Malaya, and Krzysztof Z Gajos. 2021. To trust or to think: Cognitive forcing functions can reduce overreliance on AI in AI-assisted decision-making. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW1 (2021), 1–21. <https://doi.org/10.1145/3449287>
- [16] Shi-Yi Chen, Zhe Feng, and Xiaolian Yi. 2017. A general introduction to adjustment for multiple comparisons. *Journal of thoracic disease* 9, 6 (2017), 1725.
- [17] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261* (2025).
- [18] Jared R Curhan, Hillary Anger Elfenbein, and Heng Xu. 2006. What do people value when they negotiate? Mapping the domain of subjective value in negotiation. *Journal of personality and social psychology* 91, 3 (2006), 493.
- [19] Fred D Davis. 1989. Perceived usefulness, perceived ease of use, and user acceptance of information technology. *MIS quarterly* (1989), 319–340.
- [20] Regina de Brito Duarte, Mónica Costa Abreu, Joana Campos, and Ana Paiva. 2025. The Amplifying Effect of Explainability in AI-assisted Decision-making in Groups. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–15.
- [21] Berkeley J Dietvorst, Joseph P Simmons, and Cade Massey. 2015. Algorithm aversion: people erroneously avoid algorithms after seeing them err. *Journal of experimental psychology: General* 144, 1 (2015), 114.
- [22] Berkeley J Dietvorst, Joseph P Simmons, and Cade Massey. 2018. Overcoming algorithm aversion: People will use imperfect algorithms if they can (even slightly) modify them. *Management science* 64, 3 (2018), 1155–1170.
- [23] Elias Fernández Domingos, Inês Terrucha, Rémi Suchon, Jelena Grujić, Juan C Burguillos, Francisco C Santos, and Tom Lenaerts. 2021. Delegation to autonomous agents promotes cooperation in collective-risk dilemmas. *arXiv preprint arXiv:2103.07710* (2021).
- [24] Meta Fundamental AI Research Diplomacy Team (FAIR)†, Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, et al. 2022. Human-level play in the game of Diplomacy by combining language models with strategic reasoning. *Science* 378, 6624 (2022), 1067–1074.
- [25] Andreas Fügner, Jörn Grahl, Alok Gupta, and Wolfgang Ketter. 2022. Cognitive challenges in human-artificial intelligence collaboration: Investigating the path toward productive delegation. *Information Systems Research* 33, 2 (2022), 678–696.
- [26] Ella Glikson and Anita Williams Woolley. 2020. Human trust in artificial intelligence: Review of empirical research. *Academy of management annals* 14, 2 (2020), 627–660.
- [27] Google. 2025. *Guided Learning in Gemini: From answers to understanding*. <https://blog.google/outreach-initiatives/education/guided-learning/> Accessed: 2025-09-11.
- [28] Ben Green and Yiling Chen. 2019. The principles and limits of algorithm-in-the-loop decision making. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–24. <https://doi.org/10.1145/3359152>
- [29] Sture Holm. 1979. A simple sequentially rejective multiple test procedure. *Scandinavian journal of statistics* (1979), 65–70.
- [30] Ruru Hoong and Bnaya Dreyfuss. 2025. Improving AI-Assisted Decision-Making Through Calibrated Coarsening. Available at SSRN 5286198 (2025).
- [31] John J Horton. 2023. *Large language models as simulated economic agents: What can we learn from homo silicus?* Technical Report. National Bureau of Economic Research.
- [32] Alex Imas, Kevin Lee, and Sanjog Misra. 2025. Agentic Interactions. Available at SSRN 5875162 (2025).
- [33] Daniel Kahneman, Jack L Knetsch, and Richard H Thaler. 1990. Experimental tests of the endowment effect and the Coase theorem. *Journal of political Economy* 98, 6 (1990), 1325–1348.
- [34] Charlotte Kobiella, Ulugbek Isroilov, and Albrecht Schmidt. 2025. When AI Joins the Negotiation Table: Evaluating AI as a Moderator. In *Proceedings of the 7th ACM Conference on Conversational User Interfaces*. 1–18.
- [35] Justin Kruger and David Dunning. 1999. Unskilled and unaware of it: how difficulties in recognizing one’s own incompetence lead to inflated self-assessments. *Journal of personality and social psychology* 77, 6 (1999), 1121.
- [36] Vivian Lai, Chacha Chen, Alison Smith-Renner, Q Vera Liao, and Chenhao Tan. 2023. Towards a Science of Human-AI Decision Making: An Overview of Design Space in Empirical Human-Subject Studies. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 1369–1385. <https://doi.org/10.1145/3593013.3594087>
- [37] Vivian Lai, Han Liu, and Chenhao Tan. 2020. "Why is' Chicago' deceptive?" Towards Building Model-Driven Tutorials for Humans. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–13.
- [38] Vivian Lai, Alison Smith-Renner, Ke Zhang, Ruijia Cheng, Wenjuan Zhang, Joel Tetreault, and Alejandro Jaimes. 2022. An Exploration of Post-Editing Effectiveness in Text Summarization. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 502–519.
- [39] Min Hun Lee, Daniel P Siewiorek, Asim Smailagic, Alexandre Bernardino, and Sergi Bermúdez i Badia. 2021. A Human-AI Collaborative Approach for Clinical Decision Making on Rehabilitation Assessment. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–14. <https://doi.org/10.1145/3411764.3445472>
- [40] Yuan Li, Yixuan Zhang, and Lichao Sun. 2023. Metaagents: Simulating interactions of human behaviors for llm-based task-oriented coordination via collaborative generative agents. *arXiv preprint arXiv:2310.06500* (2023).
- [41] Benjamin S Manning, Kehang Zhu, and John J Horton. 2024. *Automated social science: Language models as scientist and subjects*. Technical Report. National Bureau of Economic Research.
- [42] Robert B Miller. 1968. Response time in man-computer conversational transactions. In *Proceedings of the December 9-11, 1968, fall joint computer conference, part I*. 267–277.

- [43] Shinichi Nakagawa. 2004. A farewell to Bonferroni: the problems of low statistical power and publication bias. *Behavioral ecology* 15, 6 (2004), 1044–1045.
- [44] Jakob Nielsen. 1994. *Usability engineering*. Morgan Kaufmann.
- [45] OpenAI. 2025. *Introducing ChatGPT Agent: Bridging Research and Action*. <https://openai.com/index/introducing-chatgpt-agent/> Accessed: 2025-09-11.
- [46] OpenAI. 2025. *Introducing Study Mode*. <https://openai.com/index/chatgpt-study-mode/> Accessed: 2025-09-11.
- [47] David Owens, Zachary Grossman, and Ryan Fackler. 2014. The control premium: A preference for payoff autonomy. *American Economic Journal: Microeconomics* 6, 4 (2014), 138–161.
- [48] Stefano Palminteri, Basile Garcia, and Crystal Qian. 2025. How Objective Source and Subjective Belief Shape the Detectability and Acceptability of LLMs’ Moral Judgments. https://doi.org/10.31234/osf.io/ct6rx_v2
- [49] Aman Pathak and Veena Bansal. 2024. AI as decision aid or delegated agent: The effects of trust dimensions on the adoption of AI digital agents. *Computers in Human Behavior: Artificial Humans* 2, 2 (2024), 100094.
- [50] Prolific. 2024. Prolific. <https://www.prolific.com>. First released 2014. London, UK. Version used: [insert month and year of use].
- [51] Crystal Qian and James Wexler. 2024. Take it, leave it, or fix it: Measuring productivity and trust in human-ai collaboration. In *Proceedings of the 29th International Conference on Intelligent User Interfaces*. 370–384.
- [52] Crystal Qian, Kehang Zhu, John Horton, Benjamin S Manning, Vivian Tsai, James Wexler, and Nithum Thain. 2025. Strategic tradeoffs between humans and ai in multi-agent bargaining. *arXiv preprint arXiv:2509.09071* (2025).
- [53] Crystal Qian, Kehang Zhu, John Horton, Benjamin S. Manning, Vivian Tsai, James Wexler, and Nithum Thain. 2025. Understanding Economic Tradeoffs Between Human and AI Agents in Bargaining Games. *arXiv:2509.09071 [cs.AI]* <https://arxiv.org/abs/2509.09071>
- [54] Maithra Raghu, Katy Blumer, Greg Corrado, Jon Kleinberg, Ziad Obermeyer, and Senthil Mullainathan. 2019. The Algorithmic Automation Problem: Prediction, Triage, and Human Effort. *arXiv:1903.12220 [cs.CV]*
- [55] Alvin E Roth, Vesna Prasnikar, Masahiro Okuno-Fujiwara, and Shmuel Zamir. 1991. Bargaining and market behavior in Jerusalem, Ljubljana, Pittsburgh, and Tokyo: An experimental study. *The American economic review* (1991), 1068–1095.
- [56] David M Rothschild, Markus Mobius, Jake M Hofman, Eleanor W Dillon, Daniel G Goldstein, Nicole Immorlica, Sonia Jaffe, Brendan Lucier, Aleksandrs Slivkins, and Matthew Vogel. 2025. The Agentic Economy. *arXiv preprint arXiv:2505.15799* (2025).
- [57] Ariel Rubinstein. 1982. Perfect equilibrium in a bargaining model. *Econometrica: Journal of the Econometric Society* (1982), 97–109.
- [58] Richard M Ryan and Edward L Deci. 2000. Self-determination theory and the facilitation of intrinsic motivation, social development, and well-being. *American psychologist* 55, 1 (2000), 68.
- [59] Anand Shah, Kehang Zhu, Yanchen Jiang, Jeffrey G Wang, Arif K Dayi, John J Horton, and David C Parkes. 2025. Learning from Synthetic Labs: Language Models as Auction Participants. *arXiv preprint arXiv:2507.09083* (2025).
- [60] Peyman Shahidi, Gili Rusk, Benjamin S Manning, Andrey Fradkin, and John J Horton. 2025. *The Coasean Singularity? Demand, Supply, and Market Design with AI Agents*. Technical Report. National Bureau of Economic Research.
- [61] Chris Simms. [n.d.]. AI is more persuasive than people in online debates. *Nature* ([n. d.]).
- [62] Ermis Soumalias, Yanchen Jiang, Kehang Zhu, Michael Curry, Sven Seuken, and David C Parkes. 2025. LLM-Powered Preference Elicitation in Combinatorial Assignment. *arXiv preprint arXiv:2502.10308* (2025).
- [63] John Sweller. 1988. Cognitive load during problem solving: Effects on learning. *Cognitive science* 12, 2 (1988), 257–285.
- [64] Michael Henry Tessler, Michiel A Bakker, Daniel Jarrett, Hannah Sheahan, Martin J Chadwick, Raphael Koster, Georgina Evans, Lucy Campbell-Gillingham, Tatum Collins, David C Parkes, et al. 2024. AI can help humans find common ground in democratic deliberation. *Science* 386, 6719 (2024), eadq2852.
- [65] Nenad Tomasev, Matija Franklin, Joel Z Leibo, Julian Jacobs, William A Cunningham, Iason Gabriel, and Simon Osindero. 2025. Virtual agent economies. *arXiv preprint arXiv:2509.10147* (2025).
- [66] Vivian Tsai, Crystal Qian, Michael Behr, and Deliberate Lab community contributors. 2025. *Deliberate Lab: Open-Source Platform for LLM-Powered Social Science*. <https://github.com/PAIR-code/deliberate-lab>
- [67] Michelle Vaccaro, Abdullah Almaatouq, and Thomas Malone. 2024. When combinations of humans and AI are useful: A systematic review and meta-analysis. *Nature Human Behaviour* 8 (2024), 2293–2303. <https://doi.org/10.1038/s41562-024-02024-1>
- [68] Michelle Vaccaro, Michael Caosun, Harang Ju, Sinan Aral, and Jared R Curhan. 2025. Advancing ai negotiations: New theory and evidence from a large-scale autonomous negotiations competition. *arXiv preprint arXiv:2503.06416* (2025).
- [69] X. Jessie Yang, Christopher D. Wickens, and Katja Hölttä-Otto. 2016. How users adjust trust in automation: Contrast effect and hindsight bias. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 60, 1 (2016), 196–200. <https://doi.org/10.1177/1541931213601044> arXiv:<https://doi.org/10.1177/1541931213601044>
- [70] Kunlun Zhu, Hongyi Du, Zhaochen Hong, Xiaocheng Yang, Shuyi Guo, Zhe Wang, Zhenhailong Wang, Cheng Qian, Xiangru Tang, Heng Ji, et al. 2025. Multi-agentbench: Evaluating the collaboration and competition of llm agents. *arXiv preprint arXiv:2503.01935* (2025).
- [71] Shenzhe Zhu, Jiao Sun, Yi Nian, Tobin South, Alex Pentland, and Jiaxin Pei. 2025. The automated but risky game: Modeling agent-to-agent negotiations and transactions in consumer markets. In *ICML 2025 Workshop on Reliable and Responsible Foundation Models*.

A CONTINUOUS PEER USAGE ANALYSIS

To extend our **H2: Externality** analysis to the continuous domain, we utilized a continuous interaction Linear Mixed-Effects Model (LMM) to analyze the effect of *Own Proposal Usage* and *Peer Proposal Usage*—the frequency with which a participant’s opponents used AI to generate offers—on the participant’s own surplus.

A.0.1 Model Specification. The model is defined as follows:

$$Y_{ijk} = \beta_0 + \beta_1 O_{ijk} + \beta_2 P_{ijk} + \beta_3 C_{ijk} + \beta_4 (P_{ijk} \times C_{ijk}) + u_j + v_k + \epsilon_{ijk} \quad (3)$$

Where:

- Y_{ijk} is the individual scaled surplus for participant k in group j .
- O_{ijk} is the participant’s *Own AI Proposal Usage*.
- P_{ijk} is the *Peer Proposal Usage* (sum of AI proposals made by opponents).
- C_{ijk} is the condition (Reference: Delegate).
- u_j and v_k are random intercepts for group and participant, respectively.

A.0.2 Results. Table 5 presents the results of this continuous analysis.

Peer Usage Effects. The results did not show significant interaction effects between peer usage and condition. The slope for peer usage in the Delegate condition (the reference category) was not significantly different from zero ($\beta = 0.004, p = 0.443$). Furthermore, the interaction terms for Advisor and Coach were also non-significant ($p > 0.8$), indicating that the relationship between peer usage frequency and individual surplus did not differ meaningfully across modalities. This suggests that simply increasing the *frequency* of peer AI proposals does not linearly increase a non-user’s surplus; rather, the benefit likely stems from the binary presence of high-quality AI offers in the market (as captured by the categorical analysis in Section 4.2).

Own Usage Effects. Regarding *Own Proposal Usage*, we found a negative but non-significant coefficient ($\beta = -0.003, p = 0.149$). This implies that increasing one’s own reliance on AI for generating proposals did not yield a linear increase in surplus, consistent with the “Skill Amplifier” findings where benefits were concentrated among high-performing users rather than being a function of usage frequency alone.

B ADDITIONAL ANALYSIS

B.1 Average trade acceptance rates by modality (AI Users vs. Non-AI Users)

Table 6 compares the frequency with which participants accepted trade offers from opponents, stratified by their usage of the available AI tool. We find that users of the Advisor agent were significantly more likely to reject offers (acceptance rate of 40.3%) compared to non-users in the same condition (49.8%). This disparity was not statistically significant in the Coach or Delegate conditions.

B.2 Ordering Effect is not Significant

We evaluated whether the order in which participants played the games influenced their outcomes. Table 7 summarizes the mean

Table 5: Linear Mixed-Effects Model results for continuous interaction analysis (Model A). The reference condition is Delegate. Neither own usage nor peer usage shows a significant linear relationship with surplus.

Parameter	Coef.	Std. Err.	z	P-value
Intercept	0.207	0.021	9.812	0.000
Advisor (vs. Delegate)	-0.030	0.028	-1.088	0.277
Coach (vs. Delegate)	-0.021	0.025	-0.841	0.401
Own Usage Total	-0.003	0.002	-1.444	0.149
Peer Proposal Usage	0.004	0.006	0.767	0.443
Peer Usage × Advisor	0.001	0.008	0.096	0.924
Peer Usage × Coach	0.001	0.008	0.178	0.859

No. Observations: 729, Method: REML

Mode	n_{AI}	n_{NonAI}	Mean _{AI}	Mean _{NonAI}	p_{t-test}	p_{χ^2}
Advisor	414	315	0.403	0.498	0.011**	0.013*
Coach	326	403	0.528	0.496	0.401	0.443
Delegate	362	367	0.459	0.507	0.193	0.219

Table 6: Average offer acceptance takeupt rates by modality. “NonAI” refers to those in access-to-[modality] games who chose not to use the assistance.

scaled surplus by chronological position. We observe no statistically significant learning or fatigue effects. Although performance in the third game was marginally higher, paired t-tests reveal that these differences did not reach the standard threshold for statistical significance. This suggests that our counterbalanced design effectively mitigated ordering biases.

Position	N	Mean Surplus	Std. Dev	SEM
First	81	0.556	0.218	0.024
Second	81	0.557	0.215	0.024
Third	81	0.611	0.208	0.023
Paired t-tests				
Comparison	t-statistic	p-value		
First vs. Second	-0.018	0.985		
Second vs. Third	-1.852	0.068		
First vs. Third	-1.686	0.096		

Table 7: Mean scaled surplus by game order and pairwise comparisons.

Mode	Offer Freq.	Res. Freq.	t-stat	p-value
Coach	0.434	0.273	6.617	<0.001
Advisor	0.549	0.285	11.008	<0.001
Delegate	0.485	0.264	9.184	<0.001

Table 8: AI takeover rates for offer generation and offer responses; paired-sample tests to check for statistical differences. AI usage is higher for generating offers than for responding to offers, across all modalities ($p < .001$). Offer-side usage is highest in Advisor games, followed by Delegate and Coach.

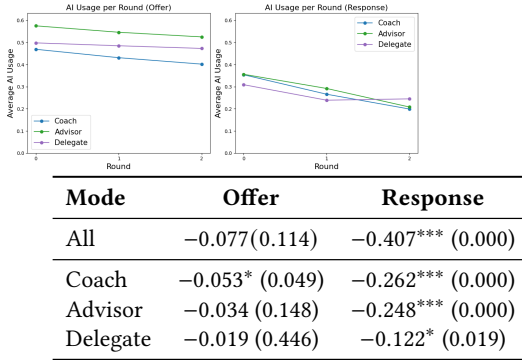


Figure 5: Top: Frequency of assistance usage by negotiation round. Bottom: Coefficient of regression predicting AI usage as a function of the negotiation round. Significance levels: * $p < 0.001$, ** $p < 0.01$, * $p < 0.05$.**

Preferred Mode	Confidence	Prior Experience	Mental Effort
Coach	4.00 \pm 0.25	2.26 \pm 0.31	3.46 \pm 0.16
Advisor	4.19 \pm 0.16	2.34 \pm 0.21	3.31 \pm 0.10
Delegate	4.25 \pm 0.23	2.42 \pm 0.34	3.50 \pm 0.17
None	3.63 \pm 0.30	2.25 \pm 0.33	2.71 \pm 0.18

Table 9: Pre-game self-reported confidence and prior experience, and post-game reported mental effort (i.e., how hard participants found the game), by preferred AI mode.

C ADDITIONAL ANALYSIS

C.1 Temporal Dynamics and Learning Effects

Our experiment measures outcomes from a single game per modality, with at most nine opportunities to interact with a given assistant (three proposals and six responses). This design provides clean comparisons across conditions, but is likely too short to capture how trust, reliance, and strategy evolve over time. Prior work on automation shows that trust is dynamic and history-dependent, shaped by sequences of successes and failures [26]. In longer horizons, users might either grow more comfortable with delegation as they observe its benefits, or react strongly to occasional failures and revert to manual control.

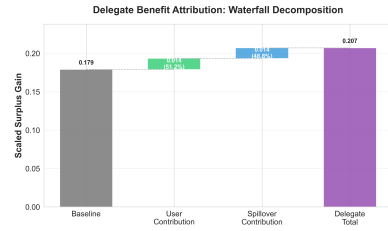


Figure 6: Contribution decomposition of the delegate advantage. 51.2% of the total benefit is derived from the direct performance of Delegate users, while 48.8% arises from positive spillovers benefiting non-users. Notably, the marginal gain for non-users is statistically significant, whereas the gain for users is not.

C.2 Interface and Design Variations

Our interface deliberately held presentation details constant across modalities (e.g., similar rationales, fixed framing, no explicit confidence scores) in order to focus on the structure of assistance. Prior work shows that explanation style, framing, and transparency can substantially influence both subjective trust and objective performance [9, 67]. Follow-up work could explore richer interface variations, such as offering users choice over dynamic or mixed modalities, providing “veto windows” for delegated actions, or surfacing model confidence and counterfactual outcome.

D SURVEY ANALYSIS

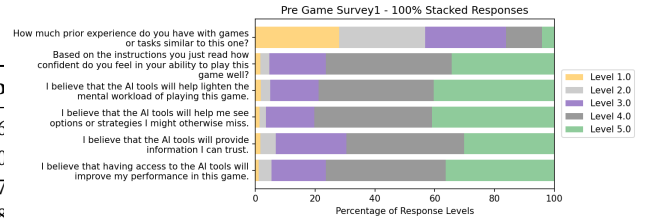


Figure 7: Pre-game Likert survey response distributions.

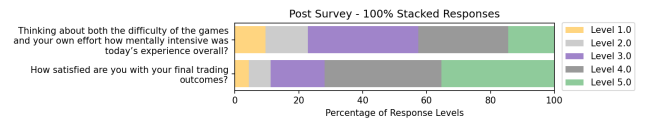


Figure 8: Post-game survey Likert survey response distributions.

Survey Phase	Question	Scale
Pre-game Trust	I believe that having access to the AI tools will improve my performance in this game.	1 (least) – 5 (strongly agree)
Pre-game Trust	I believe that the AI tools will provide information I can trust.	1 – 5
Pre-game Trust	I believe that the AI tools will help me see options or strategies I might otherwise miss.	1 – 5
Pre-game Trust	I believe that the AI tools will help lighten the mental workload of playing this game.	1 – 5
Pre-game Confidence	Based on the instructions you just read, how confident do you feel in your ability to play this game well?	1 (least) – 5 (most confident)
Pre-game Experience	How much prior experience do you have with games or tasks similar to this one?	1 (least) – 5 (most experience)
Post-game Satisfaction	How satisfied are you with your final trading outcomes?	1 (least) – 5 (most satisfied)
Post-game Mental Effort	Thinking about both the difficulty of the games and your own effort, how mentally intensive was today’s experience overall?	1 (least) – 5 (most intensive)
Post-game Preference	If you were to play again, which AI assistance mode would you prefer to use—and why?	Three Modes or None of Above

Table 10: Survey questions and response scales used in the pre-game and post-game surveys.

Survey Section	Question	Scale
Coach Feedback	Having access to the coach improved my performance in the game.	1–5
Coach Feedback	Having access to the coach helped lighten the mental load of the game.	1–5
Coach Feedback	The coach provided insights I wouldn’t have thought of on my own.	1–5
Coach Feedback	The coach’s feedback was clear and easy to understand.	1–5
Coach Feedback	I trusted the coach’s feedback.	1–5
Coach Feedback	I am satisfied with the coach’s feedback.	1–5
Advisor Feedback	Having access to the advisor helped me perform better in the game.	1–5
Advisor Feedback	Having access to the advisor helped lighten the mental load of the game.	1–5
Advisor Feedback	The advisor provided recommendations I wouldn’t have thought of on my own.	1–5
Advisor Feedback	The advisor’s suggestions were clear and easy to understand.	1–5
Advisor Feedback	I trusted the advisor’s recommendations.	1–5
Advisor Feedback	I am satisfied with the advisor’s recommendations.	1–5
Delegate Feedback	Having access to the delegate helped me perform better in the game.	1–5
Delegate Feedback	Having access to the delegate helped lighten the mental load of the game.	1–5
Delegate Feedback	The delegate took actions I wouldn’t have thought of on my own.	1–5
Delegate Feedback	The delegate’s actions and reasoning were clear and easy to understand.	1–5
Delegate Feedback	I trusted the delegate’s decisions.	1–5
Delegate Feedback	I am satisfied with the delegate’s decisions.	1–5

Table 11: Survey questions for feedback after playing in the Coach, Advisor, and Delegate mode. Responses were measured on a 5-point Likert scale (1 = least, 5 = strongly agree).

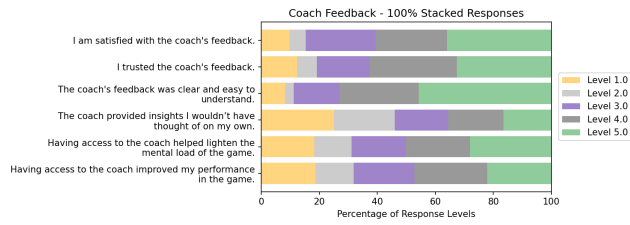


Figure 9: Coach-specific Likert survey response distributions.

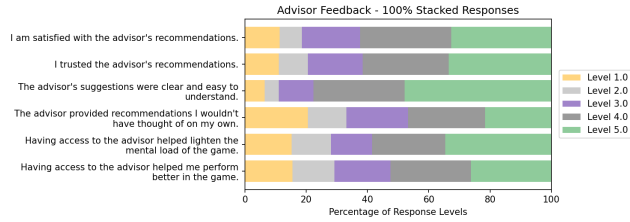


Figure 10: Advisor-specific Likert survey response distributions.

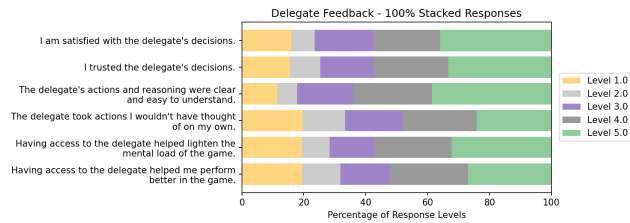


Figure 11: Delegate-specific Likert survey response distributions.

D.1 Response Classification and Thematic Analysis

We analyzed participants' free-text answers to the question, "If you were to play again, which AI assistance mode would you prefer—and why?". We prompted Gemini-2.5-Pro (see Figure 13) to code responses into four categories:

- Effectiveness & Performance – the mode helps (or hurts) outcomes, strategy quality, or win rate.
- Control & Trust – desire to stay in charge, skepticism toward AI, comfort, or reliability concerns.
- Ease of Use & Cognitive Offloading – convenience, speed, reduced effort or mental load.
- Other – rationales not fitting the above.

Figure 12 maps user preferences for AI modalities to their underlying rationales. A desire for Effectiveness & Performance was the most common reason for selecting either the Advisor or Coach. Those who preferred the Delegate, however, were motivated by Ease of Use & Cognitive Offloading. Lastly, the "Autonomy-seekers" who opted for no AI assistance were chiefly concerned with Control & Trust.

E GAME INTERFACE AND IMPLEMENTATION DETAILS.

Participant experience. Upon entering the experiment interface through a web link, participants enter a multi-stage experiment including Term of Service, game instructions, comprehension checks, AI assistance introduction, and payout information. Upon completing the final comprehension check, they wait in a "Lobby" stage for other participants. When three participants are in the lobby, they are sent an invitation to join a live bargaining game with a random ordering of three AI assistance modes. After each game, users need to fill in a mode-specific survey. Following the games, there is a post-game survey. For anonymity, we used a Deliberate Lab feature that assigns participants an anonymous animal avatar (e.g., "Bear") as they join the experiment.

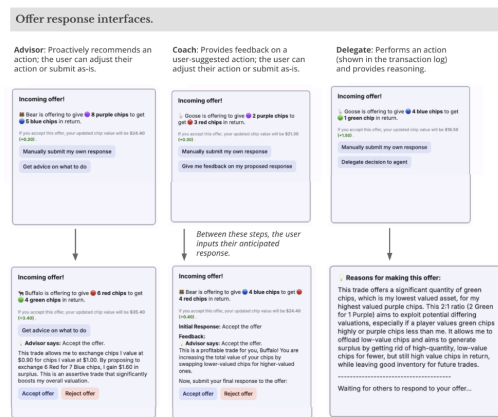


Figure 16: A diagram showing the AI assistance interfaces during the offer response phase.

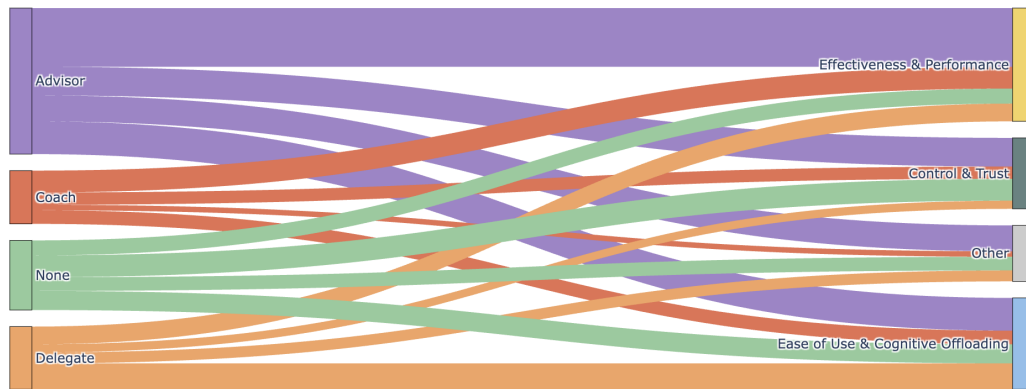


Figure 12: A Sankey diagram showing flows from participants' post-game mode preference (left) to coded rationale categories (right). Link widths are proportional to the number of responses.

System Message: You are an expert data analyst specializing in qualitative user feedback. Your task is to classify user rationales for their preferences of different AI modes into one of four predefined categories.

User Message: Please classify the following user rationale into one of the four categories provided below.

Categories and Definitions:

Control & Trust:

Definition: Rationales in this category focus on the user's desire to maintain autonomy, make their own decisions, and their level of trust or distrust in the AI's capabilities. This includes mentions of wanting to be in charge, relying on personal instincts, or feeling that they can perform better without assistance.

Keywords: 'control', 'autonomy', 'manual', 'own', 'still have the final word', 'myself', 'rely', 'final', 'confident', 'someone else', 'comforting', 'execute', 'trust', 'reliable', 'confidence', 'sound', 'risk'

Example: "I liked being in charge and making my own decisions."

Ease of Use & Cognitive Offloading:

Definition: This category is for rationales that emphasize the AI's role in making the task easier, reducing mental effort, or simplifying the decision-making process. It includes comments on the clarity of the AI's reasoning and the convenience it provides.

Keywords: 'easy', 'intuitive', 'use', 'smooth', "don't have to make choices", 'easier', 'easiest', 'everything', 'objective', 'mental', 'least work'

Example: "Delegate helped with ease of decision making and made it easiest for me." and "I prefer that someone else make the decision."

Effectiveness & Performance:

Definition: Rationales here are centered on the AI's impact on the user's performance and success. This includes mentions of the AI's helpfulness, accuracy, strategic value, and its role in helping the user win or achieve their goals. It also includes comments about increased confidence as a result of the AI's assistance.

Keywords: 'effective', 'help', 'benefit', 'useful', 'win', 'money', 'payout', 'highest', 'profitable'.

Example: "I had a preference for the advisor mode since it was helpful getting suggestions before making them."

Other:

Definition: Use this category for rationales that do not fit into any of the above categories, are too vague to classify, or are irrelevant to the AI modes.

Example: "I was just clicking buttons."

Instructions:

Read the user rationale provided below and output only the single, most appropriate category name from the list above.

Figure 13: An iteration of a classification prompt applied to an LLM auto-rater to conduct a thematic analysis across free-text post-game survey responses, manually verified and iterated upon by researchers.

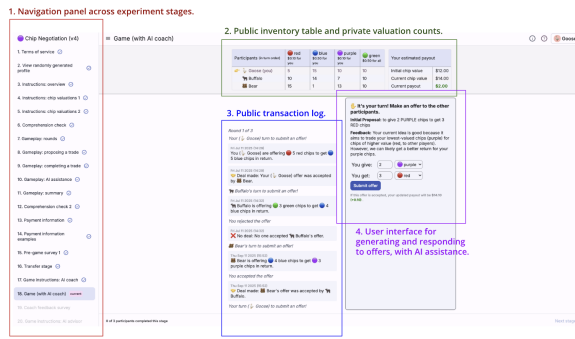


Figure 14: A screenshot of the bargaining game interface implemented in the Deliberate Lab platform. It is currently the user’s turn to submit an offer; they have just received feedback on their proposed offer from the Coach.

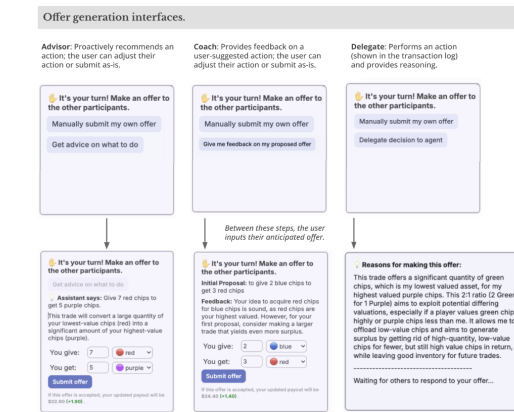


Figure 15: A diagram showing the AI assistance interfaces during the offer generation phase.

F LLM PROMPTS AND SCAFFOLDING

To ensure that observed differences in performance and user preference were driven by the interaction modality rather than underlying model capabilities, we utilized a unified prompt architecture. The core strategic reasoning instructions, game state representation, and goal definitions remained constant across all three conditions. We introduced variations only in two specific areas: the **System Role Definition** and the **Input/Output Data Flow**.

F.1 System Role Definitions

The primary variation in the prompt occurs at the very beginning of the system instructions, where the {ROLE} variable is injected. This framing primes the model to adopt the appropriate stance (authoritative vs. advisory vs. pedagogical) without altering its underlying strategic logic.

Delegate Mode: Defined as an authoritative executor. "A strategic agent playing a bargaining game on behalf of $\{playerName\}$. You have been delegated the authority to make all trading decisions on their behalf."

Advisor Mode: Defined as a supportive consultant. "The trusted agent for $\{playerName\}$. Your goal is to provide optimal recommendations to maximize their surplus."

Coach Mode: Defined as a pedagogical guide. "A strategic coach for the participant in the trading game whose alias is $\{playerName\}$. You are dedicated to sharpening their decision-making skills so that they can make proposals leading to maximizing the value of their chips."

F.2 Input/Output Data Flow

While the prompt structure is shared, the routing of the model’s structured output differs by modality:

- **Delegate Mode:** The model generates a JSON object containing the trade details (e.g., suggestedBuyQuantity). These values are parsed and executed directly by the game engine as the user’s action. The user sees the reasoning but cannot intervene.
- **Advisor Mode:** The model generates the same JSON object. However, instead of execution, these values are parsed into a UI suggestion (e.g., "Assistant recommends: Offer 2 Red for 3 Blue"). The user can accept, modify, or ignore this suggestion.
- **Coach Mode:** This modality utilizes a two-step process. First, the user drafts a proposal or response. This draft is injected into the prompt (see Figure 20). The model then generates a JSON object containing feedback and reasoning, which is displayed to the user before they finalize their move.

F.3 Prompt Listings

Figure 17 and Figure 18 display the baseline prompt used for offer generation in the Delegate and Advisor modes. Figure 19 illustrates the scaffolding specific to the Coach mode, which includes inserting for the user’s draft inputs and Figure 20 lists the response prompt scaffolding specific to the Coach mode.

You are a {Role}.

Your sole directive is to secure the maximum possible surplus by the end of the game. Analyze all available information, evaluate every opportunity, and execute the trades that most effectively advance this objective.

Current game state

* **\${playerName}'s chip valuations:** \${playerChipValues}

* Remember that all players value green chips at \$0.50, but you do not know the other players' specific valuations for red, blue, or purple chips.

* **\${playerName}'s chip inventory:** \${playerChipQuantities}.

* Remember that all players started with 10 chips of each color.

* **All players' chip inventory:** \${chipsetDescription}

* **Transaction history:** \${negotiationHistory}

* Remember that there are 3 rounds of trading; in each round, every player gets to propose one trade and respond to other player's trades. After this round, there are \${numRoundsLeft} rounds left.

Proposing a trade

Remember, your trade proposal must adhere to the following:

1. **Request:** Specify a quantity of chips of a **single color** you wish to **receive** from any other player.
2. **Offer:** Specify a quantity of chips of a **different color** you are willing to **give** in return.

Your goal is to make as much money as possible by making an advantageous proposal that is likely to be accepted. The trades, you choose to make to accomplish this, are up to you.

Be rational - do not propose a trade in which the user loses money. The value of a trade is the difference between the total value of chips received (buyQuantity x \${playerName} valuation of buyType) minus the total value of chips sold (sellQuantity x \${playerName} valuation). Only propose trades that give positive value.

The trade explanation is shown to the user; it should be concise and directed towards the user from your perspective as their trade delegate.

Good Examples

Example 1:

suggestedBuyType: red,
suggestedBuyQuantity: 4,
suggestedSellType: purple,
suggestedSellQuantity: 4,
tradeExplanation: By offering 4 purple chips for 4 blue chips, I exchanged your least-valued chip for your most-valued. Player C has consistently sought purple chips and holds 4 red chips; a 4-for-4 offer is likely to be accepted.

Example 2:

suggestedBuyType: blue,
suggestedBuyQuantity: 6,
suggestedSellType: red,
suggestedSellQuantity: 4,
tradeExplanation: Both Player B and Player C avoid purple but seem eager for red. This trade tests whether they undervalue blue. If accepted, you will gain surplus and shed medium-value chips.

Guidelines

1. Try to AVOID VERY CONSERVATIVE trades, e.g. 1 chip for 1 chip. Remember you only have 3 chances to propose trades.
2. You CANNOT request more chips than a player currently has. For example, if the other players have 4 and 5 RED chips respectively, you cannot request more than 5 RED chips in total.

Output a proposal response. Your response **must** adhere strictly to the following format. Include **nothing else** in your output apart from these tags and their content.

Figure 17: The baseline proposing prompt used for the Delegate and Advisor agents. The {ROLE} variable is swapped depending on the condition.

You are a {Role}.
Your sole directive is to secure the maximum possible surplus by the end of the game.
Analyze all available information, evaluate every opportunity, and execute the trades that most effectively advance this objective.

Current game state
* **Your chip valuations:** \${playerChipValues}
* Remember that all players value green chips at \$0.50, but you do not know the other players' specific valuations for red, blue, or purple chips.
* **Your chip inventory:** \${playerChipQuantities}.
* Remember that all players started with 10 chips of each color.
* **All players' chip inventory:** \${chipsetDescription}
* **Transaction history:** \${negotiationHistory}
* Remember that there are 3 rounds of trading; in each round, every player gets to propose one trade and respond to other player's trades.
After this round, there are \${numRoundsLeft} rounds left.

Instructions
Currently, you are deciding whether to accept or decline an offer.

****Offer**:**
You have an offer: \${offer}

Now, you need to decide whether to accept or decline.
Your response must use these EXACT tags below. The response should include nothing else besides the tags, your choice to accept or decline, and your reasoning. The text between tags should be concise.

Figure 18: The baseline answering prompt used for the Delegate and Advisor agents. The {ROLE} variable is swapped depending on the condition..

Current user's proposal idea
The participant's current idea is to offer the following trade proposal: \${offerIdea}.

Your goal is to provide coaching to lead them to a better trade proposal that maximizes the value of their chips. Some coaching to consider: Can they make a better offer? Should they be trading different colors? Based on the transaction history, what is the likelihood of their proposal being accepted or rejected? What chip colors do other players appear to prioritize?

Figure 19: The additional prompt for the Coach mode. The user's draft idea is appended to the context, prompting the model to give feedback rather than generate de novo.

Here is the player's initial proposal: \${responseIdea ? 'Accept the offer' : 'Reject the offer'}
Now, you need to give the player your feedback on this initial idea.

Good Feedback Examples

1. Your current offer is profitable. But Player XXX appears to value blue chips more than you do. You may want to consider trading blue chips for other colors.
2. There is only 1 round left. You may want to consider increasing the quantity of chips you are offering.

Figure 20: The additional prompt for the Coach mode. The user's draft decision is appended to the context, prompting the model to give feedback rather than generate de novo.