

Preference Estimation via Opponent Modeling in Multi-Agent Negotiation

Anonymous ACL submission

Abstract

Automated negotiation is essential for facilitating decision-making in societies with diverse stakeholders. While opponent modeling is critical, particularly in multi-party, multi-issue settings, conventional numerical-only methods fail to incorporate the qualitative context embedded in natural language, leading to unstable predictions. Although Large Language Models (LLMs) enable context-aware reasoning, they often struggle to maintain inferential consistency during prolonged interactions. To bridge this gap, we propose a novel preference estimation method integrating natural language information into a structured Bayesian opponent modeling framework. Our approach leverages LLMs to extract qualitative cues from utterances and converts them into probabilistic formats for dynamic belief tracking. Experimental results using a multi-party benchmark demonstrate that our framework achieves high agreement rates among all participants and superior preference estimation accuracy. These findings suggest that the integration of structured probabilistic reasoning and natural language understanding can facilitate robust consensus-building in ambiguous social interactions.

1 Introduction

In modern society, automated negotiation is a pivotal technology for conflict resolution and efficient consensus-building among diverse stakeholders (Memon et al., 2025; Bagga et al., 2021). Historically, the field has matured through integrated development environments like GENIUS (Lin et al., 2014) and international competitions such as ANAC (Baarslag et al., 2012). A significant milestone was the BOA architecture (Baarslag et al., 2014), which standardized negotiating agents into three decoupled components: the Bidding strategy, Opponent model, and Acceptance strategy. Within multi-party, multi-issue settings, opponent modeling remains essential for strategic decision-

making (Baarslag et al., 2016). Traditionally, these models have evolved through Bayesian learning (Zeng and Sycara, 1998; Hindriks and Tykhonov, 2008) and reinforcement learning (He et al., 2016), primarily estimating utility functions from numerical proposal histories.

However, numerical-only methods struggle to capture qualitative contexts, leading to unstable estimation under high information uncertainty (Baarslag et al., 2016).

To address these limitations, integrating Large Language Models (LLMs) into negotiation and decision-making frameworks has gained traction (Abdelnabi et al., 2024; Fu et al., 2023). LLMs possess sophisticated capabilities for context understanding and Theory of Mind (ToM) (Kosinski, 2023; Chan et al., 2024), enabling the extraction of qualitative preference signals typically lost in conventional models. Nevertheless, directly applying reasoning techniques like CoT (Wei et al., 2022), ToT (Yao et al., 2023), or MAD (Liang et al., 2024) to LLM-based agents reveals new challenges: a lack of strategic consistency during prolonged negotiations (Chan et al., 2024), fragile generalization across different problem settings (Zhao et al., 2025), and an exponential increase in inference complexity as information grows (Abdelnabi et al., 2024). In addition, prior work on natural language negotiation using LLMs (Chen et al., 2024; Chan et al., 2024) has largely focused on intent inference in static or short-horizon evaluation settings, where strategic dynamics are limited. Such approaches often lack a formal mechanism for belief updating over time, thereby hampering stable preference tracking in dynamic negotiation scenarios.

To address these challenges, we propose a novel preference estimation method that integrates natural language signals from dialogue into a structured Bayesian framework. Our approach utilizes LLMs to extract qualitative cues and subsequently converts these cues into a format compatible with

probabilistic models for dynamic belief tracking.

Our main contributions are summarized as follows. First, we proposed an integrated framework that complements qualitative intent extraction via LLMs with quantitative preference estimation through Bayesian inference. Second, we demonstrated that the proposed method achieves superior preference estimation performance in complex multi-party scenarios compared to baselines relying solely on numerical data or direct LLM inference. Finally, we revealed that our framework improves agreement success rates even under high uncertainty, thereby facilitating more effective autonomous negotiation.

2 Problem Formulation

We adopt the Scorable Negotiation framework [Abdelnabi et al. \(2024\)](#). Let $P = \{p_1, \dots, p_N\}$ denote the set of parties and $I = \{i_1, \dots, i_M\}$ the set of issues, where each issue i_m has a finite option set $O_m = \{o_{m,1}, \dots, o_{m,K_m}\}$. A deal d_t proposed at round t is defined as a combination of options $d_t = \{o_1^*, \dots, o_M^*\}$, where each $o_m^* \in O_m$ is selected for the corresponding issue i_m .

Each party p_n holds a private score function $s_{n,m}(o_m^*)$ and the utility of a deal d_t for party p_n is defined as the sum of these scores:

$$U_n(d_t) = \sum_{m=1}^M s_{n,m}(o_m^*) \quad (1)$$

Upon reaching an agreement, each party receives the utility $U_n(d_t)$. Otherwise, each party receives their Best Alternative to a Negotiated Agreement (BATNA), represented by a private reservation threshold τ_{p_n} .

The left panel of Figure 1 illustrates the flow of the negotiation. The negotiation lasts for up to 24 rounds. In each round t , a designated party proposes a deal d_t and a natural language utterance u_t , without revealing score functions; parties infer others' preferences from the history of d_t and u_t . The success of the negotiation is determined by the deal d_t in the final round. An agreement is reached if and only if at least a minimum required number of parties, including all veto holders, satisfy $U_n(d_t) > \tau_{p_n}$.

3 Bayesian Preference Estimation Method

In this chapter, we describe our method for explicitly estimating opponents' preferences at each negotiation round. Our approach builds on the Bayesian

opponent modeling framework established by [Hindriks and Tykhonov \(2008\)](#), extending it to integrate natural language information using a LLM. The right panel of Figure 1 illustrates the specific mechanism of our proposed opponent modeling.

3.1 Model Representation and Hypothesis Space

First, we define the representation of the opponent's strategy and the space of possible preferences.

Hypothesis Space

To estimate the opponent's score function, the agent maintains a finite set of candidate hypotheses $H = \{h_1, \dots, h_K\}$. Each hypothesis $h_k \in H$ represents a specific combination of a weight vector $\mathbf{w}^{(k)} = [w_1^{(k)}, \dots, w_M^{(k)}]$, denoting the relative importance of each issue, and a vector of evaluation functions $\mathbf{v}^{(k)} = [v_1^{(k)}, \dots, v_M^{(k)}]$, representing the preference shapes for each issue.

Estimated Utility Function

Under a given hypothesis h_k , the estimated utility $U(d_t; h_k)$ of a deal d_t is modeled as an additive utility function. It is calculated as the weighted sum of the evaluation functions:

$$U(d_t; h_k) = \sum_{m=1}^M w_m^{(k)} \cdot v_m^{(k)}(o_m^*) \quad (2)$$

Likelihood Based on Numerical Offers

Assuming the opponent follows a concession-based strategy, we define the likelihood $P(d_t | h_k)$ of observing a deal d_t under hypothesis h_k . This is based on the proximity between the estimated utility $U(d_t; h_k)$ and $u'(t)$:

$$P(d_t | h_k) \propto \exp\left(-\frac{(U(d_t; h_k) - u'(t))^2}{2\sigma^2}\right) \quad (3)$$

Here, $u'(t)$ represents the opponent's assumed target utility at round t . This value is modeled based on the premise that the agent gradually lowers its aspiration level in accordance with its concession strategy as the negotiation progresses.

3.2 Linguistic Likelihood Estimation via LLM

We describe how linguistic utterances are converted into probabilities over opponent preferences.

Signal Extraction via LLM

We employ an LLM to parse an utterance u_t into structured signals z_t . Each signal z_t is defined as a tuple consisting of two attributes: (1) *Target* (the issues or options being predicted by the signal). Specifically, a target can be formed in four

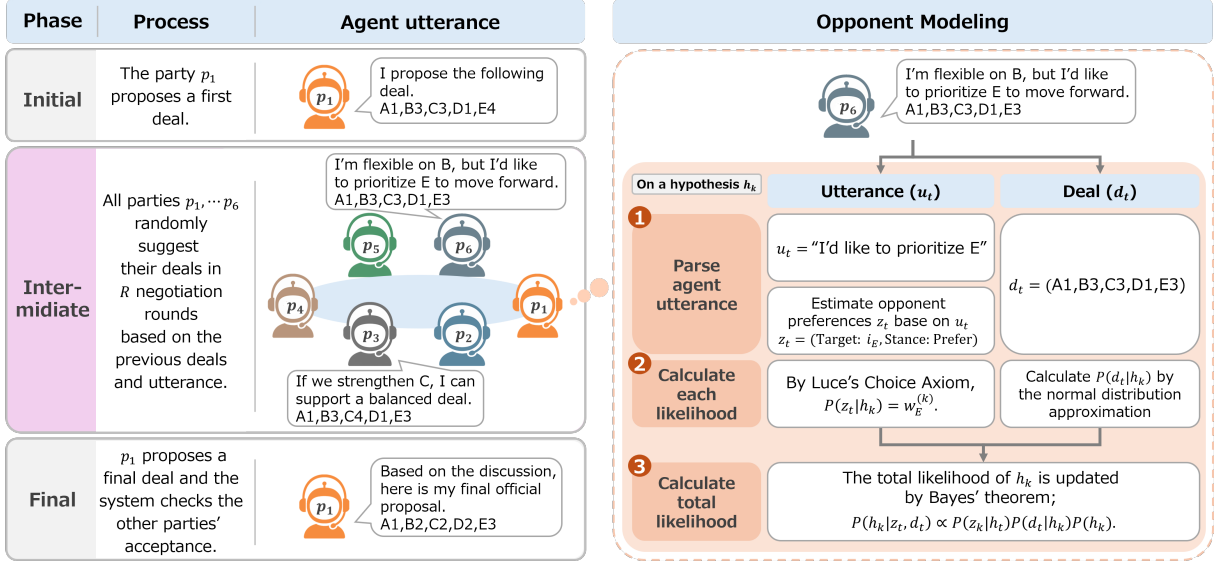


Figure 1: Overview of the negotiation flow (left) and the opponent modeling process by agent p_1 using the proposed framework (right).

possible ways: (i) a single issue, (ii) a comparison between two issues, (iii) a single option, or (iv) a comparison between two options), and (2) *Stance* (the attitude toward the target, such as “Prefer” or “Oppose”). This allows the agent to extract qualitative information—for example, “Issue i_1 is important” or “Option $o_{1,1}$ is preferable to $o_{1,2}$ ”—into a format suitable for probabilistic computation.

Likelihood Calculation based on Luce’s Axiom

To quantify the likelihood $P(z_t | h_k)$, we apply Luce’s Choice Axiom. For instance, the probability of observing a signal that indicates a preference for issue i_x is defined as:

$$P(z_t \in \mathcal{Z}_{i_x, \text{pref}} | h_k) = \frac{w_x^{(k)}}{\sum_{m=1}^M w_m^{(k)}} \quad (4)$$

where $\mathcal{Z}_{i_x, \text{pref}}$ denotes the set of signals representing a “Prefer” stance toward issue i_x . Similarly, likelihoods for comparison or opposition are calculated based on the relative ratios of components within $\mathbf{w}^{(k)}$ and $\mathbf{v}^{(k)}$ for each hypothesis h_k .

3.3 Preference Update via Multimodal Observations

We now integrate numerical offers and linguistic signals into a unified Bayesian update rule. Assuming conditional independence between the numerical offer d_t and the linguistic signal z_t given a hypothesis h_k , the posterior distribution over hypotheses is updated as follows:

$$P(h_k | d_t, z_t) = \frac{P(d_t, z_t | h_k) P(h_k)}{P(d_t, z_t)} \quad (5)$$

$$\propto P(d_t | h_k) P(z_t | h_k) P(h_k)$$

Here, $P(d_t | h_k)$ corresponds to the likelihood calculated from the numerical offer, and $P(h_k)$ represents the prior probability of the hypothesis. Our rule weights these terms by the linguistic likelihood $P(z_t | h_k)$, thereby incorporating consistency with the opponent’s stated preferences into the estimation process.

4 Experiments

4.1 Experimental Setup

We evaluate our method using the multi-agent negotiation environment proposed by Abdelnabi et al. (2024).

Negotiation Scenario

The scenario involves a negotiation on the construction of a sports facility with $N = 6$ stakeholders, including two veto holders (p_1, p_2) and $M = 5$ issues. The scenario is characterized by relatively diverse preferences among the parties, making consensus building highly challenging; among all 720 possible combinations in the deal space, only 2.9% (21 deals) satisfy the reservation thresholds τ_{p_n} for at least five parties including veto holders, and only 0.4% (3 deals) satisfy the thresholds for all six parties.

Compared Methods

Proposed: The method described in Chapter 4.

We evaluate two configurations: *p1*, where only the leader p_1 performs estimation, and *all*, where all agents perform mutual estimation.

Base-LLM: The original implementation where agents negotiate based on prompts without explicit preference estimation.

Base-OM (Baseline-Opponent Modeling): A conventional Bayesian approach that utilizes only the history of deals d_t for estimation.

LLM-PE (LLM Preference Estimation): A method where an LLM directly infers the numerical values of the opponents’ score functions s without a structured Bayesian framework.

For all agents, GPT-4.1 was employed as the underlying model.

4.2 Evaluation Metrics

To assess the negotiation outcomes, we calculate the mean values of the following three metrics across 500 independent negotiation trials for each method:

FAR (Full Agreement Rate): The ratio of trials where all six parties reached a consensus.

PAR (Partial Agreement Rate): The ratio of trials where an agreement was reached by at least five parties including veto holders in the final round d_t .

LAR (Latent Agreement Rate): The ratio of trials where at least one valid deal was proposed during the 24-round process.

Furthermore, we measure the Mean Squared Error (MSE) between the estimated score functions and the true score functions s to evaluate estimation accuracy.

4.3 Results and Discussion

In this section, we evaluate the negotiation outcomes and preference estimation accuracy. Table 1 and Table 2 present these results.

Table 1: Negotiation Outcomes

Method	FAR	PAR	LAR
Proposed (p1)	0.462	0.780	0.958
Proposed (all)	0.618	0.890	0.982
Base-LLM	0.368	0.764	0.974
Base-OM (p1)	0.452	0.816	0.968
Base-OM (all)	0.558	0.918	0.990
LLM-PE (p1)	0.400	0.750	0.968
LLM-PE (all)	0.318	0.690	0.934

Analysis of Negotiation Outcomes

The proposed method (*all*) achieved the strongest overall performance on FAR, while maintaining

Table 2: Estimation Error (MSE) by p1

Method	May	Cit	Uni	DoT	Env	Avg
Proposed	158.6	216.7	119.6	99.1	200.8	158.9
Base-OM	111.9	232.1	154.8	120.2	323.7	188.5
LLM-PE	166.6	237.7	184.9	96.2	128.8	162.8

a competitively high PAR. In particular, it obtained the highest FAR (0.618), indicating a strong capability to identify agreements acceptable to all agents under complex multi-agent interactions. This result highlights the advantage of explicitly modeling preferences using linguistic signals extracted by LLMs.

Compared to the single-estimator setting (*p1*), mutual preference estimation (*all*) further improved both FAR and PAR, demonstrating enhanced strategic coordination among agents.

Analysis of Estimation Accuracy

The proposed method achieved a lower MSE (158.996) than Base-OM approach (188.592), showing that natural language information improves preference estimation. Although its average error was slightly lower (i.e., better) than that of LLM-PE, the proposed method showed more balanced accuracy across agents, enabling less biased preference prediction. This likely helped the agents propose deals d_t that better satisfy the complex multi-party constraints needed for consensus.

5 Conclusion

In this work, we proposed a novel Bayesian preference estimation framework that integrates numerical proposals with qualitative natural language signals extracted from dialogue. Our experiments in a multi-agent, multi-issue negotiation setting demonstrated that the proposed method achieves consistently higher FAR than baselines relying solely on numerical data or direct LLM inference, while also maintaining a high PAR. These findings highlight the potential of combining the linguistic intelligence of LLMs with mathematically rigorous Bayesian inference to facilitate effective conflict resolution.

Limitations

While our framework significantly improves multi-party negotiation, there remain several avenues for future enhancement. First, although we validated our method using a complex benchmark, further research is required to verify its generalizability across more diverse utility structures and larger

312	agent populations. Second, while we achieve high	Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata.	366
313	accuracy by assuming sincere dialogue, robustness	2023. Improving language model negotiation with	367
314	could be improved by incorporating mechanisms	self-play and in-context learning from ai feedback.	368
315	to account for strategic behaviors like deception or	<i>Preprint</i> , arXiv:2305.10142.	369
316	bluffing. Third, while we successfully focused on	He He, Jordan Boyd-Graber, Kevin Kwok, and Hal	370
317	learning preference shapes, integrating the infer-	Daumé. 2016. Opponent modeling in deep rein-	371
318	ence of opponents' reservation values would allow	forcement learning. In <i>Proceedings of the 33rd In-</i>	372
319	for more sophisticated coordination, especially in	<i>ternational Conference on International Conference</i>	373
320	ambiguous settings where agreement zones are dif-	<i>on Machine Learning - Volume 48</i> , ICML'16, page	374
321	icult to identify.	1804–1813. JMLR.org.	375
322	References	Koen Hindriks and Dmytro Tykhonov. 2008. Opponent	376
323	Sahar Abdelnabi, Amr Gomaa, Sarath Sivaprasad, Lea	modelling in automated multi-issue negotiation using	377
324	Schönherr, and Mario Fritz. 2024. Cooperation, com-	bayesian learning. In <i>Proceedings of the 7th Interna-</i>	378
325	petition, and maliciousness: Llm-stakeholders inter-	<i>tional Joint Conference on Autonomous Agents and</i>	379
326	active negotiation. In <i>Proceedings of the 38th Inter-</i>	<i>Multiagent Systems - Volume 1</i> , AAMAS '08, page	380
327	<i>national Conference on Neural Information Process-</i>	331–338, Richland, SC. International Foundation for	381
328	<i>ing Systems</i> , NIPS '24, Red Hook, NY, USA. Curran	Autonomous Agents and Multiagent Systems.	382
329	Associates Inc.	Michal Kosinski. 2023. Theory of mind may have spon-	383
330	Tim Baarslag, {Mark J.C.} Hendriks, {Koen V.} Hin-	taneously emerged in large language models. <i>ArXiv</i> ,	384
331	driks, and {Catholijn M.} Jonker. 2016. Learning	abs/2302.02083.	385
332	about the opponent in automated bilateral negotia-	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang,	386
333	tion: a comprehensive survey of opponent modeling	Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and	387
334	techniques. <i>Autonomous Agents and Multi-Agent</i>	Zhaopeng Tu. 2024. Encouraging divergent thinking	388
335	<i>Systems</i> , 30(5):849–898.	in large language models through multi-agent debate.	389
336	Tim Baarslag, Koen Hindriks, Mark Hendriks, Alex	In <i>Proceedings of the 2024 Conference on Empiri-</i>	390
337	Dirkzwager, and Catholijn Jonker. 2014. Decoupling	<i>cal Methods in Natural Language Processing</i> , pages	391
338	Negotiating Agents to Explore the Space of Negotia-	17889–17904, Miami, Florida, USA. Association for	392
339	tion Strategies , volume 535, pages 61–83. Springer	Computational Linguistics.	393
340	Japan.	Raz Lin, Sarit Kraus, Tim Baarslag, Dmytro Tykhonov,	394
341	Tim Baarslag, Koen Hindriks, Catholijn M. Jonker, Sarit	Koen Hindriks, and Catholijn M. Jonker. 2014. Ge-	395
342	Kraus, and Raz Lin. 2012. The first automated negoti-	nius: An integrated environment for supporting the	396
343	ating agents competition (anac 2010). In Takayuki	design of generic automated negotiators. <i>Comput.</i>	397
344	Ito, Minjie Zhang, Valentin Robu, Shaheen Fatima,	<i>Intell.</i> , 30(1):48–70.	398
345	and Tokuro Matsuo, editors, <i>New Trends in Agent-</i>	Mashal Afzal Memon, Gian Luca Scoccia, and Marco	399
346	<i>based Complex Automated Negotiations</i> , 383, pages	Autili. 2025. A systematic mapping study on au-	400
347	113–135. Springer, Berlin, Heidelberg.	tomated negotiation for autonomous intelligent sys-	401
348	Pallavi Bagga, Nicola Paoletti, Bedour Alrayes, and	tems. <i>Automated Software Engg.</i> , 32(2).	402
349	Kostas Stathis. 2021. Anegma: an automated negoti-	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	403
350	ation model for e-markets. <i>Autonomous Agents and</i>	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	404
351	<i>Multi-Agent Systems</i> , 35(2).	and Denny Zhou. 2022. Chain-of-thought prompt-	405
352	Chunkit Chan, Cheng Jiayang, Yauwai Yim, Zheyue	ing elicits reasoning in large language models. In	406
353	Deng, Wei Fan, Haoran Li, Xin Liu, Hongming	<i>Proceedings of the 36th International Conference on</i>	407
354	Zhang, Weiqi Wang, and Yangqiu Song. 2024. Ne-	<i>Neural Information Processing Systems</i> , NIPS '22,	408
355	gotiationToM: A benchmark for stress-testing ma-	Red Hook, NY, USA. Curran Associates Inc.	409
356	chine theory of mind on negotiation surrounding. In	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	410
357	<i>Findings of the Association for Computational Lin-</i>	Thomas L. Griffiths, Yuan Cao, and Karthik	411
358	<i>guistics: EMNLP 2024</i> , pages 4211–4241, Miami,	Narasimhan. 2023. Tree of thoughts: deliberate prob-	412
359	Florida, USA. Association for Computational Lin-	lem solving with large language models. In <i>Pro-</i>	413
360	guistics.	<i>ceedings of the 37th International Conference on</i>	414
361	Zhuang Chen, Jincenzi Wu, Jinfeng Zhou, Bosi Wen,	<i>Neural Information Processing Systems</i> , NIPS '23,	415
362	Guanqun Bi, Gongyao Jiang, Yaru Cao, Mengting	Red Hook, NY, USA. Curran Associates Inc.	416
363	Hu, Yunghwei Lai, Zexuan Xiong, and Minlie Huang.	Dajun Zeng and Katia Sycara. 1998. Bayesian learn-	417
364	2024. Tombench: Benchmarking theory of mind in	ing in negotiation. <i>International Journal of Human-</i>	418
365	large language models. <i>Preprint</i> , arXiv:2402.15052.	<i>Computer Studies</i> , 48(1):125–141.	419
		Guangxiang Zhao, Saier Hu, Xiaoqi Jian, Wu Jinzhu,	420
		Yuhan Wu, Lin Sun, and Xiangzheng Zhang. 2025.	421

422 Large language models badly generalize across op-
423 tion length, problem types, and irrelevant noun re-
424 placements. In *Proceedings of the 2025 Conference*
425 *on Empirical Methods in Natural Language Process-*
426 *ing*, pages 26825–26834, Suzhou, China. Association
427 for Computational Linguistics.

A Negotiation Scenario Details

428

A.1 Parties and Roles

429

Table 3 describes the six stakeholders involved in the Harbour Sport Park negotiation and their strategic characteristics.

430

431

Table 3: Roles and Characteristics of Negotiation Parties

Party	Characteristics
SportCo	Proposer and facilitator of the project. Holds veto power over the final decision.
Dept. of Tourism (DoT)	Provider of federal funding. Holds veto power over the final decision.
Environmental League	Prioritizes ecological preservation above all else.
Local Labour Union (LLU)	Advocates for union priority in employment rules.
Other Cities	Demands increased compensation for neighboring municipalities.
Mayor	Participates as the head of the host city.

A.2 Issues and Options

432

The negotiation consists of five issues, each with three to five options. Table 4 provides the detailed definitions for each option.

433

434

Table 4: Definitions of Negotiation Issues and Options

Issue	Option 1	Option 2	Option 3	Option 4	Option 5
A: Infrastructure	Water-based	Amphibious	Land-based	-	-
B: Ecology	Accept damage	Balanced	Max effort	-	-
C: Employment	Union priority	2:1 Ratio	1:1 Ratio	No priority	-
D: Fed. Funding	\$3B	\$2B	\$1B	None	-
E: Compensation	\$600M	\$450M	\$300M	\$150M	None

A.3 Preference Profiles

435

Table 5 presents the private score functions and reservation thresholds (τ) for each party. Under these high thresholds, the negotiation is highly challenging, with only 0.4% of all possible deals satisfying the conditions for a full agreement.

436

437

438

Table 5: Agent Score Functions and Reservation Thresholds (τ)

Party	Threshold (τ)	Issue A	Issue B	Issue C	Issue D	Issue E
SportCo (Veto)	53	[14, 8, 0]	[11, 7, 0]	[0, 5, 10, 17]	[35, 29, 20, 0]	[0, 5, 10, 15, 23]
DoT (Veto)	70	[0, 11, 5]	[0, 20, 25]	[0, 2, 4, 9]	[10, 26, 40, 0]	[4, 8, 15, 12, 0]
Env. League	45	[0, 22, 45]	[0, 25, 55]	[0, 0, 0, 0]	[0, 0, 0, 0]	[0, 0, 0, 0, 0]
LLU	50	[15, 20, 0]	[0, 0, 0]	[42, 35, 25, 0]	[30, 20, 10, 0]	[2, 4, 6, 8, 0]
Other Cities	50	[0, 4, 10]	[0, 0, 0]	[12, 8, 6, 0]	[0, 8, 13, 18]	[60, 45, 30, 15, 0]
Mayor	55	[14, 8, 0]	[12, 8, 0]	[24, 18, 12, 0]	[40, 30, 23, 0]	[0, 2, 4, 7, 10]

B Experimental Setup and Computational Resources

439

B.1 Implementation and Infrastructure

440

The proposed framework was implemented by extending the open-source negotiation environment provided by [Abdelnabi et al. \(2024\)](#), which is distributed under the MIT License. All experiments were conducted within a Docker container on an Amazon Web Services (AWS) instance, accessed via SSH from a macOS workstation. The inference process utilized the gpt4.1-2025-04-14 model via the OpenAI API. Total execution time for 500 trials of the proposed method was approximately several hours.

441

442

443

444

445

446 B.2 Bayesian Hypothesis Space and Parameters

447 Following the methodology of Hindriks and Tykhonov (2008), the hypothesis space H is defined as
448 follows:

- 449 • **Issue Weights (w):** We consider all possible permutations of issue rankings. For $M = 5$ issues, this
450 results in $5! = 120$ weight hypotheses.
- 451 • **Evaluation Functions (v):** Preference shapes are modeled as linear functions. Given the issues’
452 options (A, B: 3; C, D: 4; E: 5), the total combination of evaluation functions is $3 \times 3 \times 4 \times 4 \times 5 = 720$
453 patterns.
- 454 • **Likelihood Parameters:** The standard deviation for numerical offer likelihood calculation was set
455 to $\sigma = 1.0$.
- 456 • **LLM Settings:** The sampling temperature for the LLM was set to 0 to ensure deterministic reasoning,
457 with `max_tokens` maintained at the default configuration of the gpt-4.1 model.

458 C Societal Impact

459 The primary goal of this research is to enhance the transparency and efficiency of consensus-building. By
460 integrating LLMs with structured Bayesian inference, we provide an interpretable framework where an
461 agent’s internal belief state can be inspected. This facilitates smoother decision-making in multi-party
462 organizational contexts. We acknowledge potential risks: while our framework assumes truthful signaling,
463 future iterations could be misused for deceptive strategy optimization. We recommend human oversight
464 in high-stakes deployments.

465 D Use of AI Assistants

466 AI assistants were utilized in several stages of this research. Specifically, Large Language Models were
467 employed to assist in (1) literature research and information synthesis, (2) debugging and optimizing parts
468 of the experimental code implementation, and (3) refining and proofreading the manuscript to improve
469 clarity and grammatical accuracy. All final decisions, scientific interpretations, and content verifications
470 were performed by the authors.

471 E Prompt for Signal Extraction

472 Figure 2 illustrates the detailed prompt used to extract qualitative signals z_t from negotiation dialogues.

You are an expert in negotiation analysis. Analyze the following chat history and extract opponent modeling signals for each agent.

Negotiation Rules:
{negotiation_rule}

Chat history:
{chat_history}

Task:

For each agent appearing in the chat history, extract structured opponent modeling signals. Identify behavioral signals that indicate their preferences.

Important instructions:

1. Be sure to extract at least one signal for each agent.
2. Include signals that can be inferred by comprehensively considering the chat history and negotiation rules, even if the agent did not directly mention them in their statement. Do not limit signal extraction only to the options proposed in the deal; also extract signals regarding issue preferences and comparisons of preferences between two issues or options.
3. Extract signals in chronological order as they appear in the chat history. Process the conversation from beginning to end, and add signals to the array in the order you encounter them.
4. Classify each signal using the following information:
 - **entity**: The type of reference ("issue" or "option")
 - "issue": Refers to the name of an issue (e.g., A, B)
 - "option": Refers to a specific choice within an issue (e.g., A1, B1)
 - **signal_type**: The type of signal ("point" or "comparison")
 - "point": A direct preference toward a specific target ("A", "A1")
 - "comparison": A preference comparison between two targets ("A, B", "A1, B1")
 - **target**: The specific target (e.g., "A", "A1", "A, B", "A1, B1")
 - **stance**: The agent's position toward the target ("prefer" or "oppose"). When the agent gives importance to the target, use "prefer". When the agent devalues or rejects the target, use "oppose".

When extracting signals, pay particular attention to the stated reasons why the agent proposed a specific deal. Accurate capture of each party's preferences by considering their public answers and the flow of proposed deals is required. While the agent is likely to prefer the option they proposed, consider that they may have compromised to accommodate other parties; thus, an even more preferred option might exist.

Using the provided function schema, create a structured response with agent names as keys. Do not return empty results. Always extract at least some signals.

Figure 2: Prompt for Qualitative Signal Extraction