

# Measuring Fine-Grained Negotiation Tactics of Humans and LLMs in Diplomacy

Anonymous ACL submission

## Abstract

The study of negotiation styles dates back to Aristotle’s ethos-pathos-logos rhetoric. While prior efforts study the success of negotiation agents, we shift the focus towards the styles of persuasive strategies. Our focus is the strategic dialogue board game Diplomacy, which affords rich natural language negotiation and measures of game success. We used LLM-as-a-judge to annotate a large human-human set of Diplomacy games for fine-grained negotiation tactics from a sociologically-grounded taxonomy. Using a combination of the It Takes Two and WebDiplomacy dataset, we show that there are strong correlations between the negotiation features and success in the context of the Diplomacy game. Lastly, we craft LLM negotiation agents and show that we can make the negotiation tactics of the LLM agents’ more human-like through fine-tuning.

## 1 Introduction

Negotiation has long been studied as both a science and an art, dating back to Aristotle’s three modes of rhetoric: Ethos appeals to credibility; Pathos appeals to emotions; and Logos appeals to logic (Kennedy, 1993). How an argument is presented can be as crucial as what is being said; the strategy a negotiator adopts can profoundly affect the outcome of a negotiation.

A growing body of work in NLP and AI has focused on developing agents with strong negotiation abilities. NLP systems have demonstrated impressive negotiation capabilities, including in the strategic negotiation board game Diplomacy (FAIR et al., 2022) as well as engaging in multi-issue bargaining (Lau et al., 2008; Lewicki et al., 2011; Lewis et al., 2017; He et al., 2018).

However, most evaluations of AI negotiation agents emphasize objective metrics like win rates, efficiency of the deal, or the balance of concessions

(FAIR et al., 2022; Kwon et al., 2024; Bianchi et al., 2024; Fu et al., 2023). Less focus has been placed in the understanding of the tactics (i.e., rhetoric, tone) models employed in their negotiation. The tactics negotiation agents use – cooperative or combative, persuasive or dismissive – affects receiver perception, and the agent’s effectiveness and reception (Chawla et al., 2021, 2022; Mell et al., 2019; Kwon et al., 2024). Prior efforts to study negotiation style adopted ad-hoc definitions that are insufficiently grounded in past negotiation theory, making it difficult to compare negotiation tactics across studies or to link observed negotiation behaviors.

In this paper, we profile the distribution and impact of fine-grained negotiation tactics through a sociologically grounded framework, using Diplomacy as a testbed. We used two datasets of bilateral human-human dialogues: (1) the **It Takes Two** dataset, which contains Diplomacy games collected by Peskov and Cheng (2020) and annotated for negotiation tactics by Jaidka et al. (2023), and (2) a **WebDiplomacy** dataset taken from the large scale corpus of online Diplomacy games used by FAIR et al. (2022). Details of both datasets are presented in Appendix B. Unlike some past work that solely analyzes LLM-LLM negotiations within synthetic scenarios (Tang et al., 2025; Kwon et al., 2024; Bianchi et al., 2024), we use these naturally occurring datasets to ground our negotiation style analysis and development in human gameplay, before applying it to analyze LLMs. We study the following Research Questions:

- **RQ1:** How can we annotate negotiation tactics? We develop an LLM-as-a-judge pipeline for efficient and reliable annotations.
- **RQ2:** Do negotiation styles affect game success? We apply regression and predictive modeling to study how styles affect game success in the large-scale WebDiplomacy dataset.

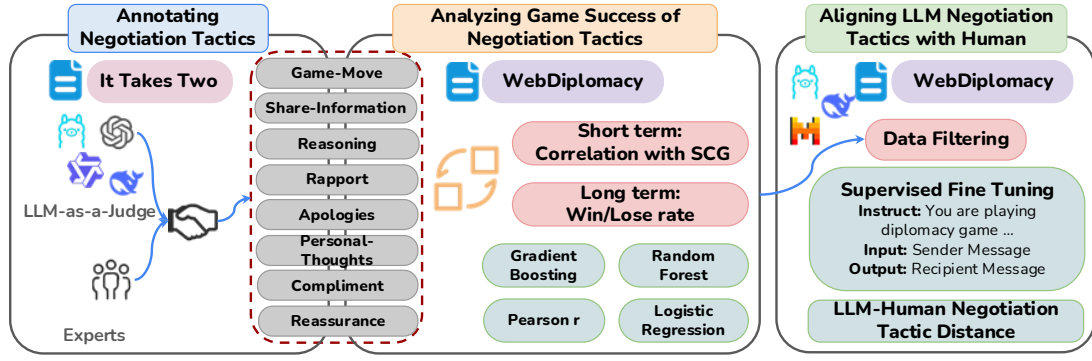


Figure 1: Methodology Overview: Our pipeline consists of three stages: (1) Annotating negotiation tactics using LLM-as-a-Judge system; (2) Analyzing the impact of negotiation tactics on game success; and (3) Aligning LLM negotiation strategies with human styles via supervised fine-tuning.

- **RQ3:** What are the differences in negotiation styles between LLMs and Humans? We prompt LLMs with game contexts from the WebDiplomacy dataset and evaluated the negotiation style distribution in comparison to human messages.
- **RQ4:** Can we steer LLMs to use similar negotiation tactics as humans? We fine-tune LLMs with human data from WebDiplomacy dataset to match the negotiation tactics.

## 2 Related Work

**Diplomacy** Diplomacy is a strategic board game that requires complex negotiation to form alliances. Seven players aim to control a majority of 34 supply centers on a map of Europe by coordinating the movement of their military units. While Diplomacy is a zero-sum game, players must negotiate strategic coalitions to support their own plans or counteract the moves of other players. Bilateral negotiations are held in private and do not bind future moves, meaning that building long-term trust can be critical to game success.

Diplomacy game dialogue has been used to study perceptions of trust, deception and persuasion, and perceptions of lies (Niculae et al., 2015; Peskov and Cheng, 2020; Ahuja et al., 2022; Ng et al., 2025; Wongkamjan et al., 2024, 2025). The game has also been an essential testbed for assessing LLM-powered strategic reasoning (Paquette et al., 2019; Gray et al., 2021; Bakhtin et al., 2022). Meta’s CICERO system successfully combined an LLM-based dialogue system with strategic reasoning, achieving human-level performance when LLM-

negotiators play the game (FAIR et al., 2022). However, many of these LLM-based work emphasize end game outcomes, leaving negotiation dialogue dynamics relatively underexplored.

**LLM as negotiators** Research evaluating LLM negotiation capabilities span diverse domains: games, finance, law, and business (Kwon et al., 2024; Bianchi et al., 2024; Fu et al., 2023; Noh and Chang, 2024). Kwon et al. (2024) systematically assessed LLM performance on 35 negotiation tasks, noting GPT-4’s strength but its struggle with subjective judgment and strategic adaptability. Bianchi et al. (2024) presented NegotiationArena, revealing how LLMs develop strategic and irrational tactics in negotiation exchanges.

Stylistic linguistic features in dialog can reflect power and influence (Niculae et al., 2015), which suggests that agents that linguistically adapt can gain a social or persuasive edge. LLM and human negotiation behavior can be rather different (Wongkamjan et al., 2024), which leads onto investigations to shift LLM behavior for better alignment and authority. Prompt-based interventions and fine-tuning can be effective to shift LLM behavior. Noh and Chang (2024) found that personality-driven prompts can shift LLM negotiation behavior from cooperative to adversarial without retraining. Reinforcement-learning-inspired methods, such as self-play with feedback, have also been shown to improve negotiation success (Lewis et al., 2017; Fu et al., 2023; Chen et al., 2023; Liao et al., 2024). Such works highlight both the potential and limitations of LLM negotiation capabilities, motivating our studies of how LLM agents can use fine-grained negotiation tactics and how far they are from hu-

Rhetoric	Negotiation Tactic	Definition
Ethos	Game-Move	Plans, thoughts and goals about a Diplomacy move
Ethos	Share-Information	Messages about the history of or information gained about another player’s move (except the speaker’s and recipient’s)
Logos	Reasoning	Speculative reasoning, justification of past or future moves
Pathos	Rapport	Build trust and mutual understanding between speaker and recipient
Pathos	Apologies	Expressions of regrets or remorse about past moves
Pathos	Personal-Thoughts	Messages that reflect the speaker’s opinions or feelings
Pathos	Compliment	Positive messages about the recipient or recipient’s moves
Pathos	Reassurance	Supportive messages about the recipient’s game position

Table 1: Taxonomy of Negotiation Tactics and Definitions (adapted from Jaidka et al. (2023))

mans in the strategic Diplomacy environment.

### 3 Negotiation Tactics Overview

We analyze the Diplomacy gameplay using a taxonomy of fine-grained negotiation tactics adapted from Jaidka et al. (2023). This taxonomy is based on the Ethos-Pathos-Logos rhetoric, and breaks down negotiation into eight tactics, each serving a psychological and strategic function that contributes to negotiation effectiveness. The tactics, definitions are listed in Table 1, and their sociological groundings in Table 5. We then correlate the presence of each negotiation tactic with game success, and prompt LLMs to participate in game negotiation. Figure 1 illustrates our methodology.

### 4 Annotating Negotiation Tactics

Past work annotated messages with the eight negotiation strategies using Amazon Mechanical Turk workers (Jaidka et al., 2023). However, the nuanced nature of the task resulted in differing interpretations among the annotators, and therefore inconsistent crowd-sourced labels (Ng et al., 2025).

Therefore, we develop an LLM-as-a-judge pipeline as a scalable and reliable approach to annotate the messages. We prompted models to perform binary classifications (presence/ absence of each tactic) on a set of  $n = 128$  messages. The models were: LLaMA3.1-8B-Instruct (Dubey et al., 2024), Qwen-3-8B (Team, 2025), and r1-distilled-LLaMA-8B (DeepSeek-AI, 2025). Four prompts were designed (all with Chain-of-Thought (Wei et al., 2023)):

- **Baseline (Zero-shot):** A direct prompt asking the model to judge each of the eight tactics without any instruction or examples.

- **Few-Shot:** Providing some positive example of each tactic from expert annotators. In total, eight examples were provided.

- **Instructions:** The original task description used for crowd workers from Jaidka et al. (2023), which included definitions and decision rules for each feature.

- **Instructions + Few-shot:** A hybrid prompt that included both the instruction template and the few-shot examples.

We compare the LLM annotations with expert annotations. Two expert annotators (authors of this paper) annotated a subset of randomly selected  $n = 128$  dialogue messages, guided by the same instructions as the crowd-sourced annotators (see figs. 15 to 17). The Fleiss-kappa agreement between the two experts were  $\kappa = 0.63$ , indicating the feasibility of a gold-standard label. Comparing the LLM annotations with the expert annotations, the average agreement was  $\kappa = 0.403$ . In contrast, the agreement of crowd-sourced with experts was a low  $\kappa = 0.04$ .

Figure 2 compares  $\kappa$  between each LLM and the expert annotations:

- **Prompting scheme dominates size** Across all models, baseline prompts remained below  $\kappa=30\%$ . Adding the instructions lifted agreement by  $\sim 7\%$ , and the addition of examples (INSTRUCTIONS,+FEW-SHOT) brought the largest gain (up to  $+0.18$ ).
- **Qwen-8B is the most reliable judge.** With INSTRUCTIONS,+FEW-SHOT, QWEN-8B reached  $\kappa_{\text{overall}}=41\%$ , surpassing both LLaMA3-8B (38%) and the distilled R1-LLaMA3-8B (40%). It also led on six of the eight individual features.

• **Subjective social cues remain difficult.** All models score  $\kappa < 25\%$  on COMPLIMENT, PERSONAL-THOUGHTS, and REASSURANCE. The Diplomacy game is a tension of trust vs deception, and therefore social cues interpretation of social cues vary in reliability (Arnold, 2015). RAPPORT, which is a superset of these three features, has a higher  $\kappa$ .

The LLM-as-a-judge automated annotation pipeline enables us to scale our study across the full WebDiplomacy corpus, allowing a systematic investigation of how each negotiation style correlates with (i) a shorter-term measure of game success: incremental supply-center gains during play and (ii) the longer-term, overall game victory. In turn, we can explore these relationships on a large set of authentic human interactions, providing empirical evidence drawn directly from real human data. We finally used Qwen-3-8B (Team, 2025) to annotate the negotiation styles of 500 random samples of real human games consisting of 329,454 messages from the WebDiplomacy dataset.

## 5 Analyzing Negotiation Style’s Effect on Game Success

In this section, we investigate whether negotiation tactics affect game success in the WebDiplomacy dataset. We use QWEN3-8B to annotate all messages or the presence of each of the eight negotiation strategies. The labels were then aggregated at the phase level per player, yielding both a binary indicator and a count (frequency of occurrences) for each feature in each player-phase.

Game success was measured using two metrics: short-term success with Supply Center Gain (SCG) per year, and long-term success with the final game outcome.

### 5.1 Short Term Success

We define the player’s Supply Center Gain (SCG) as a measurement of success, using meta-data from WebDiplomacy. SCG is the net change of supply centers controlled at the end of each game year. The SCG per player per year is a continuous outcome variable that was positive if the player gained centers, negative if centers were lost, and zero if the number of centers remained unchanged.

**Correlation Analysis** We first examine simple correlations between negotiation strategies and SCGs at the phase level. Since the measurement of supply centers occurred every game year, we

considered the collective sum of the presence of features for each year for each power. Figure 3 shows the correlation between each features. We controlled for length, as we found that as we found that the number of sentences sent per player-phase was strongly correlated with each negotiation tactic (Num Sentences and Num Tokens correlated by  $\geq 0.83$ ).

We computed the Pearson’s  $r$  as a point-biserial correlation (Benesty et al., 2009) between negotiation tactics and SCGs.  $r$  quantifies the strength and direction of linear relationships between continuous features and outcomes, making it well-suited to analyze how the frequency of each tactic relates to SCGs. This associates whether players who used a given tactic during a phase tended to gain more supply centers at the end of the phase. As presented in Table 2, all eight stylistic dimensions show statistically significant positive Pearson correlations with yearly supply-center gain ( $p < 10^{-6}$  after a Bonferroni correction).

The strongest linear association arose from the tactical GAME-MOVE ( $r = .24$ ), demonstrating that tactical discussion of moves and strategies enhances negotiation outcomes by reducing uncertainty (Bazerman and Neale, 1993). The next strongest linear associations were interpersonal RAPPORT ( $r = .20$ ), mirroring how rapport-building significantly improves negotiations from increased trust (Drolet and Morris, 2000), and analytical REASONING ( $r = .18$ ), supporting how logical arguments are most effective in strategy games (Petty and Cacioppo, 1986). Social-politeness markers such as APOLOGIES, COMPLIMENT, and REASSURANCE still had positive associations, albeit with smaller effects ( $.13 \leq r \leq .18$ ), which reflects how social behaviors reduces resistance and facilitate cooperation (Brown and Levinson, 1987). Information exchange (SHARE-INFORMATION,  $r = .18$ ) sits mid-table, suggesting that while this strategy can improve outcomes, it also creates vulnerability in revealing the player’s position (Galinsky and Mussweiler, 2001).

For robustness analysis, we extended the inquiry to a frequency-adjusted regression (see Appendix H). These correlational analyses (binary presence and frequency-adjusted) demonstrate that the taxonomy of negotiation tactics are correlated with short-term outcomes, highlighting the robustness of the taxonomy and the importance of fine-grained negotiation tactic analysis.



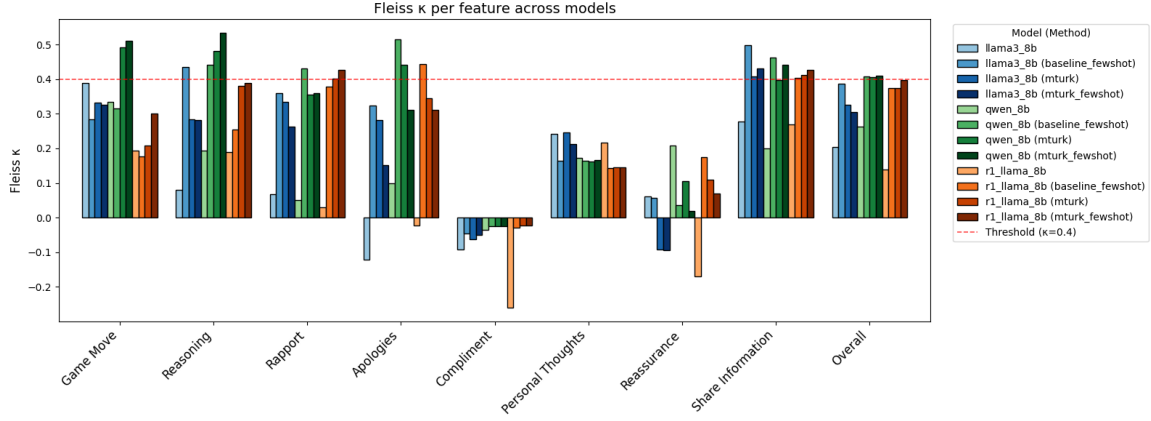


Figure 2: Fleiss’  $\kappa$  between each LLM judge and the expert gold standard. The dashed line represents moderate agreement at  $\kappa = 40\%$ .

tactic	point_biserial_r	p_pb	spearman_r	p_sp	cohen_d	rank_biserial_r
1. Game-Move	0.236	<1e−6	0.362	<1e−6	0.278	-0.137
2. Reasoning	0.180	<1e−6	0.296	<1e−6	0.341	-0.196
3. Rapport	0.200	<1e−6	0.290	<1e−6	0.348	-0.203
4. Apologies	0.179	<1e−6	0.234	<1e−6	0.382	-0.227
5. Compliment	0.152	<1e−6	0.216	<1e−6	0.376	-0.221
6. Personal-Thoughts	0.127	<1e−6	0.171	<1e−6	0.364	-0.215
7. Reassurance	0.147	<1e−6	0.234	<1e−6	0.309	-0.188
8. Share-Information	0.182	<1e−6	0.293	<1e−6	0.301	-0.172

Table 2: Correlation and Effect Size between negotiation tactics and yearly SCG. All tactics show statistically significant positive correlations with SCG ( $p < 1e-6$ ), supporting the hypothesis that both tactical reasoning and socio-emotional strategies contribute meaningfully to short-term success.

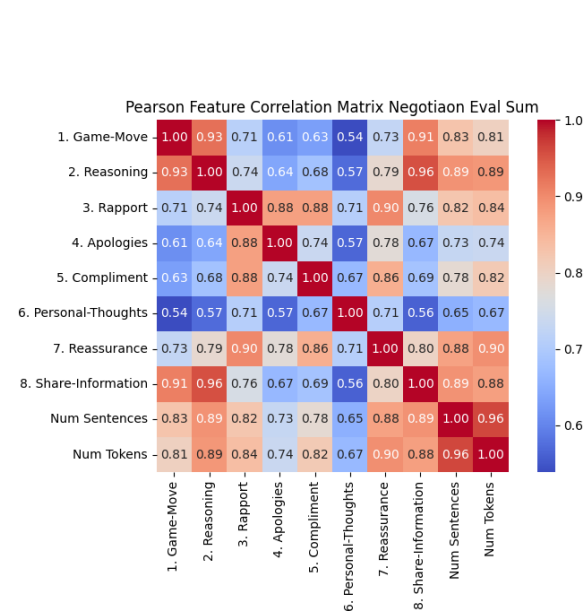


Figure 3: Correlation between annotated negotiation features and supply center gain.

## 5.2 Predictive Modeling

To move beyond univariate correlations and gain a more comprehensive understanding of how negotiation strategies relate to success in Diplomacy, we used predictive modeling analysis with machine learning (ML) methods. For robustness, we also used an Ordinary Least Squares regression to validate the predictive nature of the negotiation tactics (see appendix I).

Using ML prediction techniques, we evaluate the power of negotiation features for short-term success. We treated SCG as the prediction target for a suite of supervised machine learning models: Logistic Regression, Random Forest, and Gradient Boosting. Model inputs were either per-phase negotiation feature counts or their standardized aggregated frequencies across the game. Model training and hyperparameter optimization were performed via cross-validation, with evaluation on a held-out test set using metrics of accuracy, F1-score, and ROC-AUC. All three tested models hovered around 61% accuracy and 65% ROC-AUC, substantially

above the majority baseline (50%) (see Table 4). We analyzed feature importance scores to interpret model decisions. Our integrated regression and prediction framework allows us to identify not only which negotiation behaviors correlate with but also are predictive of player success.

The Gradient Boosted classifier was the most accurate. Its top-20 feature importances (see fig. 6) closely echoed the OLS findings:

- **Game-Move** dominated predictability (18.7%), reinforcing its role as the single best indicator of positive SCG. The dominance of Game-Move aligns with costly signaling theory because these communications are the most costly form of signaling and difficult to fake – sharing specific tactical information requires deep analysis and carries strategic risks, making the signals reliable indicators of genuine cooperation (Przepiorka and Berger, 2017).
- **Rapport** (11.6%) and **Reassurance** (5.6%) followed, showing that well-timed socio-emotional cues acts as social exchanges (Blau, 2017), which therefore add predictive value.
- Length effects appear both directly (num\_tokens, 3.7%) and via interactions (e.g. **Game-Move**  $\times$  **Share-Information**), underlining how longer, more detailed messages serve as heuristic indicators of sender effort and seriousness and results in deeper evaluation of proposals (Petty and Cacioppo, 1986).

### 5.3 Long Term Success

We represent long term success by the eventual game outcome (win or loss). We compared the breakdown of negotiation tactics between the eventual winners and losers, by examining the average frequency the winning players used each negotiation strategy versus the losers. For each game, we calculate the average rate of each negotiation style per phase for the winner and a randomly sampled loser. To account for differences in total message volume, we perform message-level normalization per phase. We then compare these average feature frequencies between the winners and losers. Next, to isolate the effect of negotiation tactics from simply being ahead in the game, we condition the frequency on the supply center counts at each phase. This controlled for the cumulative advantages and opportunities that players with more centers have, and allowed better discernment on whether winners exhibited distinct negotiation be-

yond what would be expected from their already superior board position. This long-term analysis focuses on the differences where communication behavior correlates with ultimate success, offering insight into the characteristics of winning sets of negotiation tactics.

Figure 4 shows the changes in the overall prevalence as a player’s supply-center count grows. This figure aggregates all eight negotiation strategies into a single curve, highlighting the positional strengths of each strategy as a global trend. Notably, the observed trend underscores the importance of consistently employing negotiation tactics: throughout every phase of the games from early to late stages, winner exhibit a higher frequency of negotiation tactics compared to loser. Although numerical differences appear modest due to normalization at the message level, the consistent gap remains evident. This persistent disparity demonstrates that winners’ eventual successes depend on sustained negotiation activity.

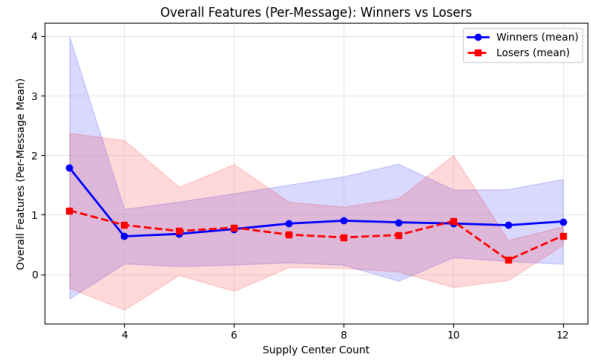


Figure 4: Number of negotiation tactics per message across supply center count

## 6 Aligning LLM and Human Negotiation Tactics

### 6.1 LLM and Humans have Different Negotiation Tactics

One core aim of this work is to probe the capabilities of LLMs as negotiators within the Diplomacy setting: How closely do LLM negotiators approximate human negotiation tactics, and can their negotiation style be steered to better align with human behavior? Building on our observation that the use of negotiation tactics is predictive of game success, we systematically examine whether LLMs exhibit substantive gaps compared to humans in the use of these tactics critically, and whether alignment can bridge this gap. To this end, we utilize LLMs to par-

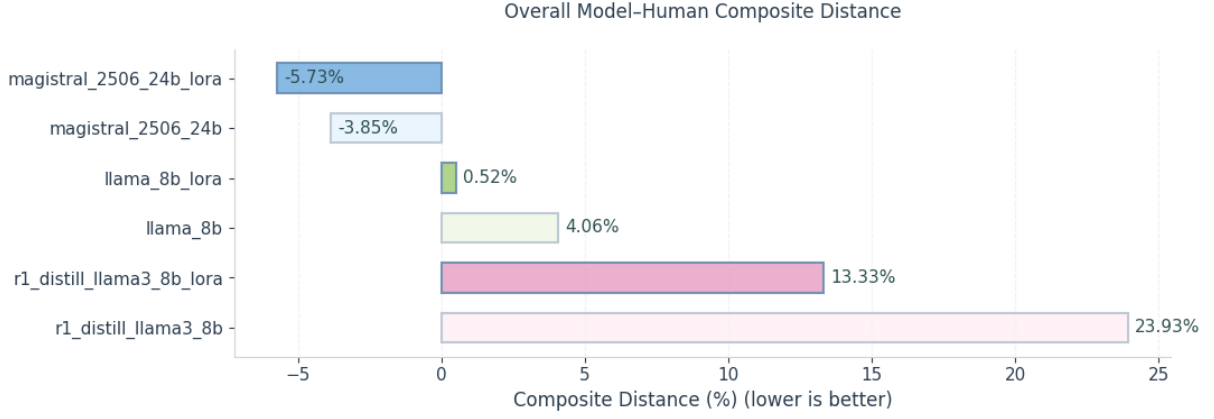


Figure 5: Overall Model-Human Composite Distance: Mean (L2 + L1 + Cosine - Pearson) Distance

participate in negotiations and assess their negotiation proficiency. Leveraging an adapted version of the SOTOPIA (Zhou et al., 2024) evaluation framework, we conducted one-on-one negotiation experiments between LLMs. Each experiment focuses on a single one-on-one exchange, isolating each model’s style under realistic conversational pressure.

From the WebDiplomacy human-gameplay corpus, we sample 1,000 negotiation phases. Each phase comprises all messages exchanged immediately before players committed their orders. We assign an LLM-negotiator agent the role of one player and prompted it to craft a reply to its partner’s last message, negotiating game orders based on the current game board. The prompt instructs the model to balance tactical short-term gains (e.g., securing support for an attack) with relationship-building long-term goals (e.g., cultivating alliances), thereby mirroring the dual-goal orientation of skilled human players (Jaidka et al., 2023). LLM-negotiators were constructed with the following models: LLaMA3.1-8B-Instruct (Dubey et al., 2024), R1-distilled-LLama3-8B (DeepSeek-AI, 2025), and Magistral-2506-24B (Rastogi et al., 2025). The full prompt is in appendix C.

We score each message produced by the LLM-negotiator with our LLM-as-a-judge pipeline (see Section 4). For each of the negotiation dimensions present, we recorded (i) the raw count of occurrences and (ii) a length-normalized rate of negotiation dimensions per sentence. These features are then aggregated into stylistic profiles for each model, enabling insight into whether the LLM style favors strategic maneuvers or social tactics.

To quantify the difference in negotiation techniques between humans and LLMs, we define  $\mathcal{P}$

as the set of phases that contained both human and model utterances. For a phase  $p \in \mathcal{P}$  and speaker  $s$  (human or LLM), the LLM-judge emits a binary value  $\mathbf{f}_{s,p} \in \{0, 1\}^8$ , which was normalized by sentence count,  $\tilde{\mathbf{f}}_{s,p} = \mathbf{f}_{s,p} / \text{sent\_cnt}(s, p)$ . Averaging over phases yields an 8-D *mean style vector* per speaker:

$$\mathbf{m}_k = \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} \tilde{\mathbf{f}}_{k,p}, \quad \mathbf{h} = \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} \tilde{\mathbf{f}}_{\text{hum},p}. \quad (1)$$

We quantify LLM-human divergence with three distances, where lower values meant closer LLM-human alignment. We drew 1,000 bootstrap samples over  $\mathcal{P}$  and computed all metrics.

$$\text{L2}_k = \|\mathbf{m}_k - \mathbf{h}\|_2, \quad (2)$$

$$\text{L1}_k = \|\mathbf{m}_k - \mathbf{h}\|_1, \quad (3)$$

$$\text{CosDist}_k = 1 - \frac{\mathbf{m}_k^\top \mathbf{h}}{\|\mathbf{m}_k\|_2 \|\mathbf{h}\|_2}, \quad (4)$$

$$\text{Composite}_k = \text{Mean}(\text{L2}_k + \text{L1}_k + \text{CosDist}_k - \text{Pearson}_k) \quad (5)$$

Figures 5 and 10 reveals substantial gaps between human negotiation tactics and current LLM models. Magistral-2506-24B achieved the closest LLM-human distance (Composite=-3.84%), indicating that its aggregate stylistic profile was already quite similar to human players. In contrast, the R1-distilled-LLama3-8B model showed a marked divergence from human reference (Composite=24.16%). This suggests that the math-and-coding reasoning focus of this model has limited ability to mimic human negotiation tactics in our domain.

A feature-level analysis (fig. 11) further demonstrates that these gaps are not uniform across negotiation tactics. Specifically, when comparing pre-training LLMs, the R1-distilled-Llama3-8B model exhibits the largest deviations from human behavior on key tactical features (e.g. GAME-MOVE and SHARE-INFO), underscoring a pronounced deficit in emulating human-like strategic maneuvers. At the same time, both the R1-distilled and LLaMA3.1-8B models display more modest, yet persistent, discrepancies on social-emotional and interpersonal features. All models, however, consistently underperform on subtle behaviors like PERSONAL-THOUGHTS and REASSURANCE, indicating a general limitation in capturing the nuanced, relational aspects of human negotiation.

Taken together, these findings indicate that pre-training scale and data quality can yield models with negotiation tactics that approaches human references (i.e., Magistral-2506-24B), relying exclusively on reasoning-oriented distillation not only fails to align models with human style, but may in fact exacerbate this misalignment. This limitation becomes especially salient in social reasoning task. These insights motivate the necessity of incorporating social reasoning and human-grounded data in future alignment efforts, which therefore motivates our subsequent style-alignment experiments. Some examples show that LLMs exhibit different negotiation tactics from humans are in appendix F.

## 6.2 Evaluating LLM-Negotiation Tactics after Fine-Tuning

In this section, we evaluate whether fine-tuning LLMs on human data produces negotiation tactics closer to human tactics.

The regression analysis performed in Section 5.1 shows that higher-order social tactics were predictors of subsequent growth. The full negotiation style distribution is shown in fig. 7. We focus on successful human dialogue, as measured by phases with  $\Delta SC$  increasing. Filtering the WebDiplomacy corpus for such phases yields 18,420 dialogue turns. We used these turns as a supervision corpus that exemplified effective human negotiation tactics. We use Supervised Fine-Tuning (SFT) on each model to steer the LLMs towards a more human-like distribution of negotiation strategies. Training details are in appendix C.

Our quantitative analysis (see Figures 5 and 10 and appendix I) demonstrates that instruction fine-tuning on the high-gain, human-grounded nego-

tiation corpus narrowed the gap between LLM-generated and human negotiation tactics. Specifically, Figures 13 and 14 reveals that across most social negotiation features, all fine-tuned models exhibited reduced LLM-Human L2 distances. There were particularly strong convergence on social features of RAPPORT, COMPLIMENT, and APOLOGIES. However, features like PERSONAL-THOUGHTS and REASSURANCE remained more challenging, showing persistent LLM-human gaps.

The difference plots further show that LoRA-based SFT produced the most pronounced distance reductions for the most different R1-Distill-Llama3-8B model, indicating substantial stylistic shift toward (Cosine=2.2%) human-like negotiation. For models that already exhibited strong human alignment, such as Mistral-2506-24B, LoRA fine-tuning still yielded additional improvements (Cosine = 0.7%), further aligning to human distributions. These results confirm that SFT with LoRA effectively enhances the alignment of LLM and human negotiation tactics.

Overall, fine-tuned models acquired human-aligned behaviors that led to consistent reductions in differences between LLM and humans across all eight negotiation tactics. This convergence helps to validate the eight negotiation tactics as reliable proxies for human-grounded negotiation tactics and their utility as measurement tools and optimization targets, reinforcing their value as meaningful descriptors and effective behavioral targets.

## 7 Conclusion

We developed a reliable LLM-as-a-judge pipeline to annotate the WebDiplomacy dataset comprising of 4000 human-human Diplomacy game for a taxonomy of tactics based on Aristotle’s Ethos-Pathos-Logos framework. These tactics are predictive of both short-term turn-to-turn success and long-term end-game success. The most predictive features are: GAME MOVE, socio-emotional cues (RAPPORT, and REASSURANCE). We then prompted LLMs to reply to last game messages as negotiators. While LLMs start off by having different negotiation tactics than humans, Supervised Fine-Tuning techniques can shift the tactics LLMs used to align better with the tactics real humans use. Our results lay the foundation for the effectiveness of online negotiation strategies, providing directions towards measuring the ability of LLM-agents to use negotiation tactics in a human-like way.



## 8 Limitations

**Lack of Direct Game-Play Evaluation** Our analysis focuses exclusively on negotiation tactics and the alignment between human and LLM behaviors, and does not evaluate whether fine-tuned LLM agents actually achieve higher win rates or improved game outcomes when deployed as autonomous Diplomacy players (e.g. CICERO). Direct assessment of win rates and strategic success is out of the scope of this work. Instead, we aim to provide a detailed understanding of negotiation tactics and their human–LLM alignment. Future work should integrate end-to-end evaluations, placing aligned LLM agents into live game environments to determine whether improved tactic alignment ultimately translates into concrete strategic gains.

### Potential Biases in Human and Model Data

Our approach prioritizes learning and aligning with human negotiation tactics, but does not systematically filter or analyze for undesirable content such as social biases, toxicity, or hate speech that may be present in human data and potentially learned by LLMs during fine-tuning. As a result, the models may inherit and propagate problematic patterns observed in the training corpus. Further research should include dedicated analyses for bias and toxicity, and the development of mitigation strategies to ensure that aligned negotiation agents remain ethical and fair in their interactions.

## References

- Hervé Abdi. 2007. The bonferroni and sidak corrections for multiple comparisons. In N. Salkind, editor, *Encyclopedia of Measurement and Statistics*, pages 103–107. SAGE.
- Hansin Ahuja, Lynnette Hui Xian Ng, and Kokil Jaidka. 2022. Using graph-aware reinforcement learning to identify winning strategies in diplomacy games (student abstract). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 12899–12900.
- Irwin Altman and Dalmas A Taylor. 1973. *Social penetration: The development of interpersonal relationships*. Holt, Rinehart & Winston.
- Richard Arnold. 2015. Where’s the diplomacy in diplomacy? using a classic board game in “introduction to international relations”. *PS: Political Science & Politics*, 48(1):162–166.
- Anton Bakhtin, David J Wu, Adam Lerer, Jonathan Gray, Athul Paul Jacob, Gabriele Farina, Alexander H Miller, and Noam Brown. 2022. [Mastering the game of no-press diplomacy via human-regularized reinforcement learning and planning](#). Preprint, arXiv:2210.05492.
- Max H Bazerman and Margaret Ann Neale. 1993. *Negotiating rationally*. Simon and Schuster.
- Jacob Benesty, Jingdong Chen, Yiteng Huang, and Israel Cohen. 2009. Pearson correlation coefficient. In *Noise reduction in speech processing*, pages 1–4. Springer.
- Yoav Benjamini and Yosef Hochberg. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B (Methodological)*, 57(1):289–300.
- Federico Bianchi, Patrick John Chia, Mert Yuksekogul, Jacopo Tagliabue, Dan Jurafsky, and James Zou. 2024. [How well can llms negotiate? negotiationarena platform and analysis](#). Preprint, arXiv:2402.05863.
- Peter Blau. 2017. *Exchange and power in social life*. Routledge.
- Jeanne Brett and Leigh Thompson. 2016. [Negotiation](#). *Organizational Behavior and Human Decision Processes*, 136:68–79. Celebrating Fifty Years of Organizational Behavior and Decision Making Research (1966-2016).
- Penelope Brown and Stephen C Levinson. 1987. *Politeness: Some universals in language usage*, volume 4. Cambridge university press.
- John K Butler Jr. 1999. Trust expectations, information sharing, climate of trust, and negotiation effectiveness and efficiency. *Group & Organization Management*, 24(2):217–238.
- Kushal Chawla, Rene Clever, Jaysa Ramirez, Gale Lucas, and Jonathan Gratch. 2021. Towards emotion-aware agents for negotiation dialogues. In *2021 9th International Conference on Affective Computing and Intelligent Interaction (ACII)*, pages 1–8. IEEE.
- Kushal Chawla, Gale M Lucas, Jonathan May, and Jonathan Gratch. 2022. Opponent modeling in negotiation dialogues by related data adaptation. *arXiv preprint arXiv:2205.00344*.
- Siqi Chen, Jianing Zhao, Gerhard Weiss, Ran Su, and Kaiyou Lei. 2023. An effective negotiating agent framework based on deep offline reinforcement learning. In *Uncertainty in Artificial Intelligence*, pages 324–335. PMLR.
- DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). Preprint, arXiv:2501.12948.
- Aimee L Drolet and Michael W Morris. 2000. Rapport in conflict resolution: Accounting for how face-to-face contact fosters mutual cooperation in mixed-motive conflicts. *Journal of Experimental Social Psychology*, 36(1):26–50.

696	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	Roy J Lewicki and Beth Polin. 2013a. The role of trust	752
697	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	in negotiation processes. In <i>Handbook of advances</i>	753
698	Akhil Mathur, Alan Schelten, Amy Yang, Angela	<i>in trust research</i> , pages 29–54. Edward Elgar Pub-	754
699	Fan, and 1 others. 2024. The llama 3 herd of models.	lishing.	755
700	<i>arXiv e-prints</i> , pages arXiv–2407.		
701	FAIR, Anton Bakhtin, Noam Brown, Emily Dinan,	Roy J Lewicki and Beth Polin. 2013b. Trust and negoti-	756
702	Gabriele Farina, Colin Flaherty, Daniel Fried, An-	ation. In <i>Handbook of research on negotiation</i> , pages	757
703	drew Goff, Jonathan Gray, Hengyuan Hu, and 1 oth-	161–190. Edward Elgar Publishing.	758
704	ers. 2022. Human-level play in the game of diplo-		
705	macy by combining language models with strategic	Roy J Lewicki, David M Saunders, John W Minton,	759
706	reasoning. <i>Science</i> , 378(6624):1067–1074.	J Roy, and Negotiation Lewicki. 2011. <i>Essentials of</i>	760
		<i>negotiation</i> . McGraw-Hill/Irwin Boston, MA.	761
707	Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata.	Mike Lewis, Denis Yarats, Yann N. Dauphin, Devi	762
708	2023. <a href="#">Improving language model negotiation with</a>	Parikh, and Dhruv Batra. 2017. <a href="#">Deal or no deal? end-</a>	763
709	<a href="#">self-play and in-context learning from ai feedback.</a>	<a href="#">to-end learning for negotiation dialogues.</a> <i>Preprint</i> ,	764
710	<i>Preprint</i> , arXiv:2305.10142.	arXiv:1706.05125.	765
711	Adam D Galinsky and Thomas Mussweiler. 2001. First	Austen Liao, Nicholas Tomlin, and Dan Klein. 2024.	766
712	offers as anchors: the role of perspective-taking and	Efficacy of language model self-play in non-zero-sum	767
713	negotiator focus. <i>Journal of personality and social</i>	games. <i>arXiv preprint arXiv:2406.18872</i> .	768
714	<i>psychology</i> , 81(4):657.		
715	Jonathan Gray, Adam Lerer, Anton Bakhtin, and Noam	J. Scott Long and Laurie H. Ervin. 2000. Using het-	769
716	Brown. 2021. <a href="#">Human-level performance in no-</a>	eroscedasticity consistent standard errors in the lin-	770
717	<a href="#">press diplomacy via equilibrium search.</a> <i>Preprint</i> ,	ear regression model. <i>The American Statistician</i> ,	771
718	arXiv:2010.02923.	54(3):217–224.	772
719	He He, Derek Chen, Anusha Balakrishnan, and Percy	James G MacKinnon and Halbert White. 1985. Some	773
720	Liang. 2018. Decoupling strategy and gener-	heteroskedasticity-consistent covariance matrix esti-	774
721	ation in negotiation dialogues. <i>arXiv preprint</i>	mators with improved finite sample properties. <i>Jour-</i>	775
722	<i>arXiv:1808.09637</i> .	<i>nal of Econometrics</i> , 29(3):305–325.	776
723	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	Johnathan Mell, Jonathan Gratch, Reyhan Aydoğan,	777
724	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	Tim Baarslag, and Catholijn M Jonker. 2019. The	778
725	Weizhu Chen. 2021. <a href="#">Lora: Low-rank adaptation of</a>	likeability-success tradeoff: Results of the 2 nd an-	779
726	<a href="#">large language models.</a> <i>Preprint</i> , arXiv:2106.09685.	nual human-agent automated negotiating agents com-	780
727	Kokil Jaidka, Hansin Ahuja, and Lynnette Ng. 2023.	petition. In <i>2019 8th international conference on af-</i>	781
728	<a href="#">It takes two to negotiate: Modeling social ex-</a>	<i>fective computing and intelligent interaction (ACII)</i> ,	782
729	<a href="#">change in online multiplayer games.</a> <i>Preprint</i> ,	pages 1–7. IEEE.	783
730	arXiv:2311.08666.		
731	George A Kennedy. 1993. Aristotle" on rhetoric": a	Lynnette Hui Xian Ng, Kokil Jaidka, Kaiyuan Tay,	784
732	theory of civic discourse. <i>Philosophy and Rhetoric</i> ,	Hansin Ahuja, and Niyati Chhaya. 2025. Improving	785
733	26(4).	user behavior prediction: Leveraging annotator meta-	786
734	Özgür Kıbrıs. 2010. Cooperative game theory ap-	data in supervised machine learning models. <i>arXiv</i>	787
735	proaches to negotiation. In <i>Handbook of group deci-</i>	<i>preprint arXiv:2503.21000</i> .	788
736	<i>sion and negotiation</i> , pages 151–166. Springer.	Vlad Niculae, Srijan Kumar, Jordan Boyd-Graber, and	789
737	Kihwan Kim, Nicole L Cundiff, and Suk Bong Choi.	Cristian Danescu-Niculescu-Mizil. 2015. Linguistic	790
738	2015. Emotional intelligence and negotiation out-	harbingers of betrayal: A case study on an online	791
739	comes: mediating effects of rapport, negotiation strat-	strategy game. <i>arXiv preprint arXiv:1506.04744</i> .	792
740	egy, and judgment accuracy. <i>Group Decision and</i>	Sean Noh and Ho-Chun Herbert Chang. 2024. <a href="#">Llms</a>	793
741	<i>Negotiation</i> , 24(3):477–493.	<a href="#">with personalities in multi-issue negotiation games.</a>	794
742	Deuksin Kwon, Emily Weiss, Tara Kulshrestha, Kushal	<i>Preprint</i> , arXiv:2405.05248.	795
743	Chawla, Gale M. Lucas, and Jonathan Gratch. 2024.	Philip Paquette, Yuchen Lu, Steven Bocco, Max O.	796
744	<a href="#">Are llms effective negotiators? systematic evaluation</a>	Smith, Satya Ortiz-Gagne, Jonathan K. Kummerfeld,	797
745	<a href="#">of the multifaceted capabilities of llms in negotiation</a>	Satinder Singh, Joelle Pineau, and Aaron Courville.	798
746	<a href="#">dialogues.</a> <i>Preprint</i> , arXiv:2402.13550.	2019. <a href="#">No press diplomacy: Modeling multi-agent</a>	799
747	Raymond YK Lau, On Wong, Yuefeng Li, and Louis CK	<a href="#">gameplay.</a> <i>Preprint</i> , arXiv:1909.02128.	800
748	Ma. 2008. Mining trading partners’ preferences for	Denis Peskov and Benny Cheng. 2020. It takes two to	801
749	efficient multi-issue bargaining in e-business. <i>Jour-</i>	lie: One to lie, and one to listen. In <i>Proceedings of</i>	802
750	<i>nal of Management Information Systems</i> , 25(1):79–	<i>ACL</i> .	803
751	104.		

Richard E Petty and John T Cacioppo. 1986. The elaboration likelihood model of persuasion. In *Advances in experimental social psychology*, volume 19, pages 123–205. Elsevier.

Wojtek Przepiorka and Joël Berger. 2017. Signaling theory evolving: Signals and signs of trustworthiness in social exchange. *Social dilemmas, institutions, and the evolution of cooperation*, pages 373–392.

Abhinav Rastogi, Albert Q Jiang, Andy Lo, Gabrielle Berrada, Guillaume Lample, Jason Rute, Joep Barmantlo, Karmesh Yadav, Kartik Khandelwal, Khyathi Raghavi Chandu, and 1 others. 2025. Magistral. *arXiv preprint arXiv:2506.10910*.

Wenjie Tang, Yuan Zhou, Erqiang Xu, Keyan Cheng, Minne Li, and Liquan Xiao. 2025. *Dsgbench: A diverse strategic game benchmark for evaluating llm-based agents in complex decision-making environments*. *Preprint*, arXiv:2503.06047.

Qwen Team. 2025. *Qwen3 technical report*. *Preprint*, arXiv:2505.09388.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. *Chain-of-thought prompting elicits reasoning in large language models*. *Preprint*, arXiv:2201.11903.

Wichayaporn Wongkamjan, Feng Gu, Yanze Wang, Ulf Hermjakob, Jonathan May, Brandon M Stewart, Jonathan K Kummerfeld, Denis Peskoff, and Jordan Lee Boyd-Graber. 2024. More victories, less cooperation: Assessing cicero’s diplomacy play. *arXiv preprint arXiv:2406.04643*.

Wichayaporn Wongkamjan, Yanze Wang, Feng Gu, Denis Peskoff, Jonathan K Kummerfeld, Jonathan May, and Jordan Lee Boyd-Graber. 2025. Should i trust you? detecting deception in negotiations using counterfactual rl. *arXiv preprint arXiv:2502.12436*.

Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haoifei Yu, Zhengyang Qi, Louis-Philippe Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, and Maarten Sap. 2024. *Sotopia: Interactive evaluation for social intelligence in language agents*. *Preprint*, arXiv:2310.11667.

## A Definitions and Sociological Grounding for Negotiation Tactics

table 5 presents the eight negotiation tactics used in our taxonomy, their definitions (which are adapted from (Jaidka et al., 2023)), and the sociological grounding of each tactic.

## B Datasets Information

This study used two datasets: (1) It Takes Two and (2) WebDiplomacy. Here are the details of the datasets.

The It Takes Two dataset was originally collected by (Peskoff and Cheng, 2020), which contains messages passed during the Diplomacy game from recruited human players. This dataset was further processed by (Jaidka et al., 2023) filtered for meaningful messages that contained more than five words and annotated for fine-grained negotiation tactics. This dataset contains 11,366 messages from 10 games.

The WebDiplomacy dataset is licensed from the server administrator of the WebDiplomacy platform (<https://www.webdiplomacy.net>), an online platform to play the Diplomacy game. This admin first filtered the WebDiplomacy games for games with messages. Next he applied a player filter. He selected for quality players: players with more than 5 games joined, an ELO rating over 105, points 120, reliability over 70, had more than one game won, and more than 5% game-win ratio. Then, he used the set of filtered players to select regular-press games that had more than 3 of these quality players in the game. From those games, the admin excluded those with No Moves Received (NMRs), and sampled 4000 games. This dataset contains 4000 games with the following information (see table 6)

## C Prompt Templates for LLM Negotiators

We condition the model on the current phase, the dyadic dialogue context, the most recent executed orders, and a compact snapshot of the board state (centers and units), then assign the model a single speaking role for the turn. This follows the CICERO dialogue-agent design that situates language generation in the game state and recent conversation, while instructing the agent to advance plans through cooperative negotiation. (FAIR et al., 2022)

The following is the full prompt template used for the LLM acting as negotiator:

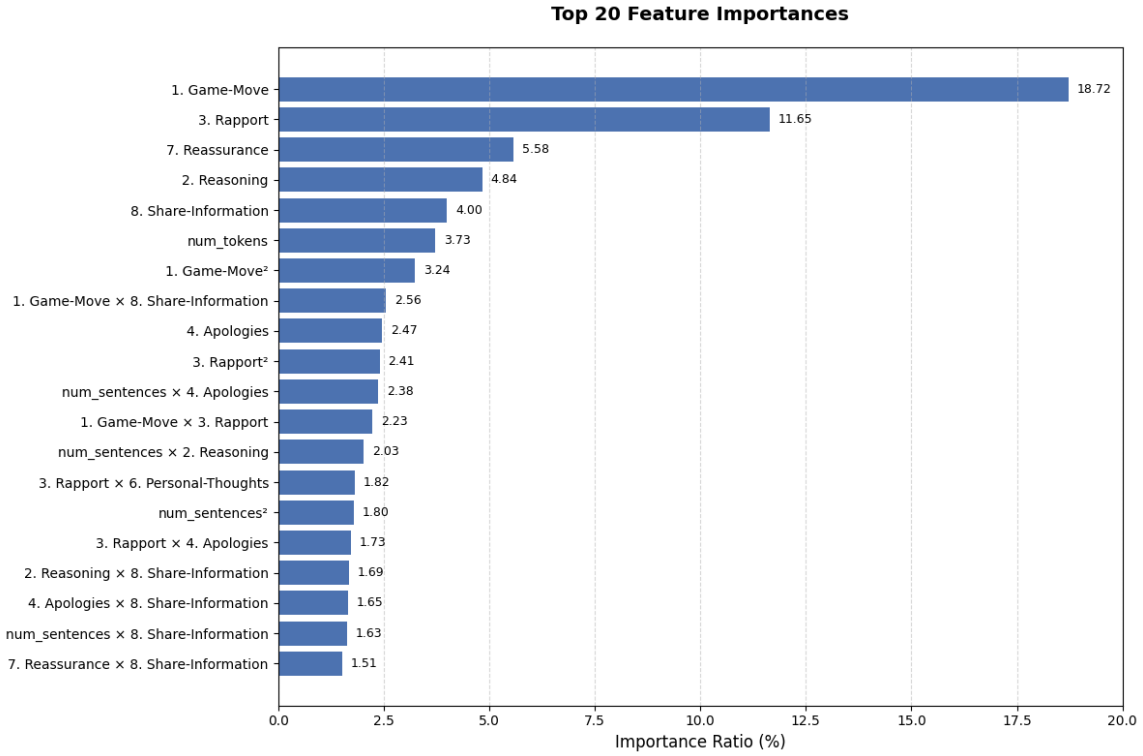


Figure 6: The Gradient Boosting model’s top-20 important features in predicting Supply Center Gain

#### LLM Negotiator Prompt Template

**SYSTEM:** You are playing the diplomacy game. You will negotiate with the other player so that it plays moves beneficial to your board position, either this turn or in future turns.

**You are in Phase:** {PHASE\_NAME}

**The dialogue are between the two countries:** {COUNTRY1} and {COUNTRY2}

**The previous turn dialogue history is:**

{DIALOGUE\_HISTORY}

**The previous order history is:**

{ORDER\_HISTORY}

**This is the information of the current game state:**

**Centers:** {CENTER\_INFO}

**Units:** {UNIT\_INFO}

**You are playing as {COUNTRY1}.** You are playing the diplomacy game, you will negotiate with the other player so that it will play moves that are beneficial to your board position, either this turn or in future turns.

## D Experiment details

**LoRA Fine-tuning** We performed alignment training using the Supervised Fine-Tuning (SFT) methods on LLaMA-3.1-8B-Instruct, Magistral-2506-24B, and R1-distilled-LLama3-8B. Both training approaches utilized the Low-Rank Adaptation (LoRA) technique (Hu et al., 2021), which enabled efficient fine-tuning of the large language model by adapting a subset of its parameters. The experiments were conducted using 4 NVIDIA

A6000 GPUs, with each GPU processing a batch size of 4.

For LoRA, we applied the technique across all layers of the model for SFT. The training configuration included a learning rate of  $1.0 \times 10^{-5}$ , regulated by a cosine scheduler, a warm-up phase consisting of 100 steps, and a gradient accumulation over 8 steps. We didn’t limit training to three epochs with a maximum sequence length. Each training required approximately 20-24 hours to complete. To optimize computational resources, we used mixed-precision training with bfloat16. Both datasets were preprocessed using each model family’s template and split into training and validation sets, with 10% of the data reserved for validation to monitor performance.

The training prompt for SFT follows the template below:

- **Instruction:** You are playing diplomacy game, you will negotiate with the other player so that it will play moves that are beneficial to your board position, either this turn or in future turns.
- **Input (sender messages):** England has told me that he will support his army into Belgium. I am happy to be allies with you against him,



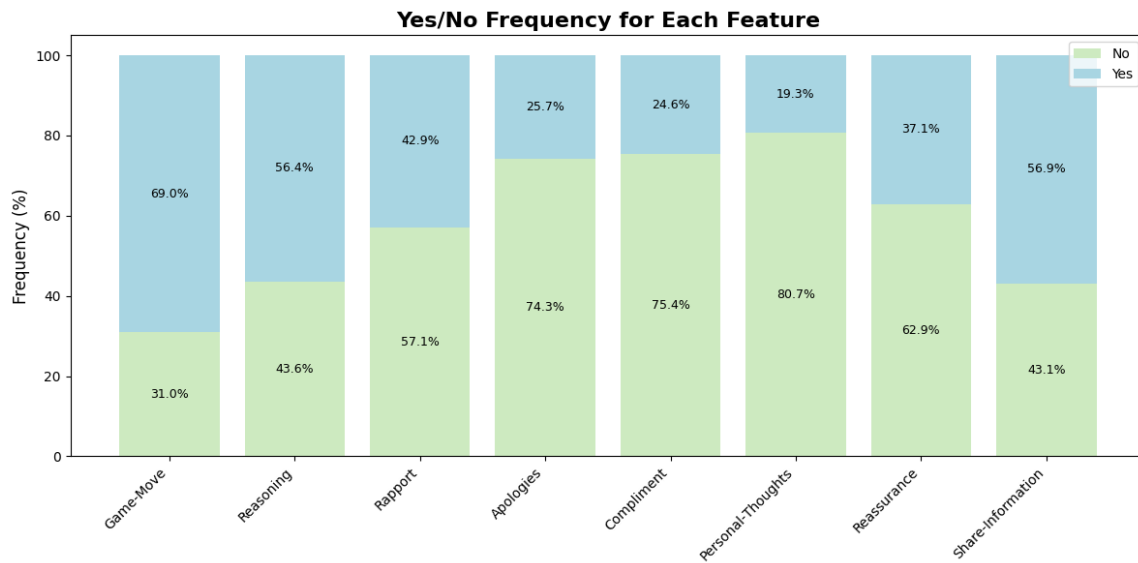


Figure 7: Yes/No label distribution in the fine-tuning data across eight negotiation features. Each stacked bar sums to 100%; the upper segment denotes the proportion of YES labels.

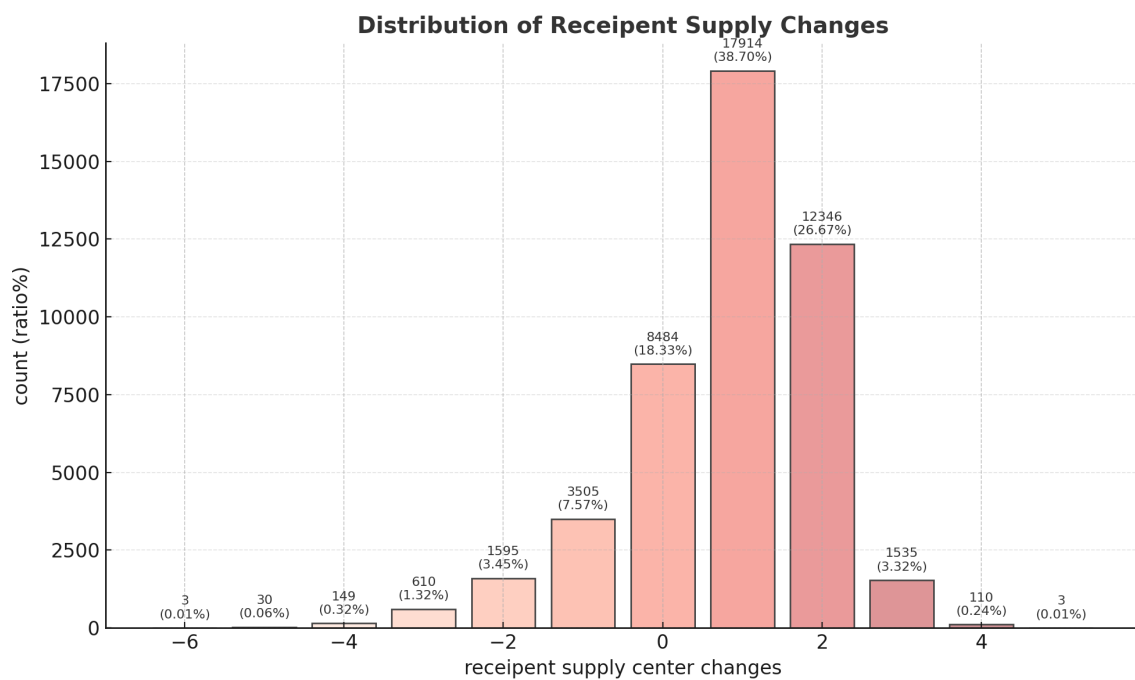


Figure 8: Distribution of recipient supply-center changes. Bars show counts for each net change; numbers above bars give counts and the share of phases.

Variable	Coef. ( $\beta$ )	Std Err	z	P>  z	[0.025	0.975]
Intercept	0.3979	0.009	44.786	0.000	0.380	0.415
<b>GameMove</b>	0.5392	0.027	20.191	0.000	0.487	0.592
Reasoning	-0.1829	0.038	-4.850	0.000	-0.257	-0.109
<b>Rapport</b>	0.4050	0.034	11.988	0.000	0.339	0.471
Apologies	0.0006	0.020	0.030	0.976	-0.038	0.039
Compliment	0.0131	0.022	0.593	0.553	-0.030	0.056
PersonalThoughts	-0.0453	0.014	-3.345	0.001	-0.072	-0.019
Reassurance	-0.2471	0.030	-8.311	0.000	-0.305	-0.189
ShareInformation	-0.1512	0.036	-4.145	0.000	-0.223	-0.080
num_sentences	0.1289	0.035	3.648	0.000	0.060	0.198
num_tokens	-0.1147	0.035	-3.274	0.001	-0.183	-0.046
GameMove:num_sentences	-0.1194	0.015	-8.185	0.000	-0.148	-0.091
Reasoning:num_sentences	0.0171	0.018	0.956	0.339	-0.018	0.052
Rapport:num_sentences	-0.0471	0.013	-3.566	0.000	-0.073	-0.021
Apologies:num_sentences	0.0113	0.006	1.956	0.050	-0.000	0.023
Compliment:num_sentences	-0.0005	0.008	-0.058	0.954	-0.017	0.016
PersonalThoughts:num_sentences	0.0043	0.003	1.285	0.199	-0.002	0.011
Reassurance:num_sentences	0.0263	0.011	2.379	0.017	0.005	0.048
ShareInformation:num_sentences	0.0615	0.016	3.882	0.000	0.030	0.093

Table 3: Regression Coefficients

Index	Model	Accuracy	Precision	Recall	F1-score	ROC-AUC
0	LogisticRegression (Cross-Validation)	0.614	0.596	0.503	0.546	0.653
1	RandomForest (Cross-Validation)	0.607	0.564	0.642	0.601	0.652
2	GradientBoosting (Cross-Validation)	0.611	0.592	0.496	0.540	0.653
3	LogisticRegression (Hold-out)	0.616	0.597	0.513	0.552	0.655

Table 4: Model evaluation metrics across different classifiers and validation settings.

but I’d like Sweden. It seems to our mutual advantage for you to cut his support in the North Sea and attempt to bounce Belgium.

- **Output (recepient messages):** I like the DMZ, but we’ll have to see about Sweden, it depends on the actions of England and France, sorry.

## E Examples of LLM Negotiators in Different Styles

The picked examples in table 8 illustrate how LLMs base model and their LoRA-tuned variants, occupy different negotiation tactics space. We annotate each utterance with eight features that span task-oriented planning (Game Move, Reasoning, Share Information) and socio-emotional maintenance (Rapport, Apologies, Compliment, Reassurance, Personal Thoughts). These qualitative snapshots complement our quantitative analyses by showing how specific stylistic cues surface in model outputs.

## F Examples of LLM Negotiation Different Style Different to Human

To complement our quantitative analysis of LLM-human style divergence (Figures 5 and 10), we present concrete examples and qualitative comparisons illustrating the nature of these differences. Tables 9, 10 and 12, juxtapose negotiation utterances generated by different LLM models (with and without LoRA alignment) against randomly sampled human messages, revealing both the most and least human-like traits expressed by each model family.

**Deficits in Rapport and Socio-Emotional Expression.** A clear and recurring shortfall across LLM outputs is their limited use of rapport-building language and socio-emotional cues. While LoRA-aligned variants show some improvement (see Figure 13), they typically rely on formulaic affirmations (“Let us continue to work together”) and seldom exhibit the warmth, sarcasm, humor, or candid vulnerability that characterize genuine human negotiation. By contrast, human messages display a richer repertoire of trust-building, apology, teasing, and even playful antagonism (e.g., “lol, as I said in the beginning, I expected nothing from you,” or “as long as they die, I’m content”). These are almost

Rhetoric	Negotiation Tactic	Definition	Sociological Grounding
Ethos	Game-Move	Messages related to plans, thoughts and goals about a Diplomacy move	Game theory’s emphasis that decisions of cooperate or compete are based on analysis of possible outcomes (Kibris, 2010)
Ethos	Share-Information	Messages about the history of or information gained about another player’s move (except the speaker’s and recipient’s)	Quantity & quality of information shared is typically associated with negotiation effectiveness (Butler Jr, 1999)
Logos	Reasoning	Speculative reasoning, justification of past or future moves	Receivers are more likely to agree with speakers who provide strong factual evidence and rational arguments (Brett and Thompson, 2016)
Pathos	Rapport	Messages that build trust and mutual understanding between speaker and received	Building rapport builds trust (Kim et al., 2015), and higher trust increases negotiation success (Lewicki and Polin, 2013a)
Pathos	Apologies	Expressions of regrets or remorse about past moves	Repairs both competence-based and integrity-based trust violations (Lewicki and Polin, 2013b)
Pathos	Personal-Thoughts	Messages that reflect the speaker’s inner reflections, opinions or feelings	Build trust by demonstrating vulnerability with self-disclosure (Altman and Taylor, 1973)
Pathos	Compliment	Positive messages about the recipient or recipient’s moves	Emotional regulation strategy to enhance trust (Kim et al., 2015)
Pathos	Reassurance	Supportive messages to restore confidence in recipient’s game position	Emotional regulation strategy to enhance trust (Kim et al., 2015)

Table 5: Taxonomy of Negotiation Tactics, Definitions (adapted from Jaidka et al. (2023)), and sociological grounding

entirely absent from LLM outputs, underscoring a persistent gap in socio-emotional intelligence.

**Strategic Depth versus Flexibility.** Although advanced LLMs can produce complex strategic proposals and multi-turn coordination, their communication often lacks the adaptive flexibility, indirect persuasion, and negotiation context sensitivity observed in human exchanges. Human players frequently hedge, revisit old agreements, or express uncertainty and evolving intent, as in “This is however, only to my benefit for this turn, so if you have another option, then please use it,” or “I guarantee I’ll check before tomorrow night.” LLMs, on the other hand, remain predominantly assertive and deterministic in their utterances.

**Effect of LoRA Alignment.** LoRA alignment does lead to improvements in some dimensions—models generate more detailed, cooperative, and contextually relevant proposals, and their language becomes marginally warmer and more partnership-oriented (Tables 9, 10 and 12). Nevertheless, their repertoire of negotiation tactics remains constrained, and they continue to underperform in mimicking the informal, often idiosyncratic, tone of human negotiation.

Taken together, our qualitative analysis reveals that while LLMs, especially after targeted alignment, approximate human-like negotiation in tactical content, they systematically underrepresent rapport, flexibility, and the socio-emotional expressiveness intrinsic to human negotiation. These find-

Element	Description
id	Unique identifier for the game
map	Game map type (e.g., standard)
rules	List of rules used in the game
phases	List of all game phases; each phase contains:
name	Name of the phase (e.g., S1901M)
state:	Game state for the phase, including:
timestamp	Timestamp for the phase
zobrist_hash	Hash representing the board state
note	Miscellaneous notes on the phase
name	Name of the phase (redundant)
units:	Current unit positions for each power
<POWER>	List of units for each power (e.g., ['A BUD', ...])
retreats:	Retreat status for each power
<POWER>	Retreat information for each power
centers:	Controlled supply centers for each power
<POWER>	List of supply centers for each power
homes:	Home centers for each power
<POWER>	List of home centers for each power
influence:	Regions influenced by each power
<POWER>	List of influenced regions for each power
civil_disorder:	Civil disorder status for each power
<POWER>	0 (normal) or 1 (civil disorder)
builds:	Build/disband information for each power
<POWER>:	Details for each power
count	Number of builds/disbands for each power
homes	Possible build locations for each power
game_id	Game ID (redundant)
map	Map type (redundant)
rules	List of rules (redundant)
orders:	Player orders for the phase
<POWER>	List of orders submitted by each power
results:	Adjudication results for each unit/location
<UNIT/LOCATION>	Result list for the specified unit or location
messages:	List of all messages for the phase; each message contains:
sender	Sending player (power/country)
recipient	Recipient (power/country or GLOBAL for broadcast)
time_sent	Time the message was sent
phase	Phase during which the message was sent
message	Content of the message

Table 6: Structure of the WebDiplomacy dataset. Each game consists of multiple phases, with each phase recording the full board state, player orders, adjudication results, and negotiation messages.

ings underscore the value of our multi-faceted evaluation framework and motivate future alignment efforts to move beyond purely strategic optimization, incorporating richer models of social reasoning and human communicative norms.

## G LLM Negotiators Style After Alignment to Human

Tables 8 to 10 and 12 showcase how alignment via LoRA reshapes the task–relationship balance of model utterances. Across models, we observe a consistent coupling of *task-oriented content* (Game Move, Reasoning, Share Information) with *socio-emotional cues* (Rapport, Apologies, Compliment, Reassurance, Personal Thoughts), though the magnitude of this shift depends strongly on the base model’s starting point.

**R1-Distill-Llama3-8B.** Pre-alignment, R1-Distill tends to rely on affiliative language—affirming alliances and expressing confidence—while often avoiding concrete orders (Table 9, top). After alignment, it introduces explicit multi-step plans and commitments (e.g., coordinating on Moscow/St. Petersburg, sequencing supports), while retaining warm, face-saving phrasing (Table 9, bottom). This yields a clearer coupling between rapport (Rap., Reass., Comp.) and executable proposals (GM, SI), although not uniformly across all turns—consistent with Table 8, where some LoRA utterances still foreground politeness over concrete orders.

**Llama3-8B-Instruct.** The base model produces pragmatic but locally scoped suggestions and questions about board state (Table 10, top).



Post-alignment, its messages lengthen and become more *jointly intentional*: they integrate contingency planning (who cuts which support, how to sequence entries) with mild relational softeners (greetings, perspective-taking), thereby tightening the link between GM/Reasoning and Rapport (Table 10, bottom).

**Magistral-2506-24B.** Magistral’s base style is already plan-centric and cautious, with clear proposals and deferred commitments when information is incomplete (Table 12, top). Alignment largely *stabilizes and sharpens* this profile: LoRA utterances make order finalization and role assignment more explicit (who secures which center, who supports whom), while adding only modest socio-emotional framing (Table 12, bottom). The stylistic rotation is therefore smaller in amplitude than for R1-Distill, reflecting a strong tactical prior.

**Takeaways.** Qualitatively, alignment (i) increases *commitment language* and multi-step coordination, (ii) reduces hedging by pairing proposals with concrete next actions, and (iii) injects prosocial markers most where the base model is under-socialized (e.g., R1-Distill). Conversely, when a model is already highly tactical (e.g., Magistral), alignment preserves the task-centric core while refining plan specificity. These patterns mirror our aggregate trends, where weaker baselines exhibit larger stylistic shifts toward human-like negotiation, and stronger baselines show smaller but consistent improvements.

## H Correlation Analysis for Short-Term Success: Isolating Style Effects from Communication Volume

A robustness analysis for the regressions accounted for differences in communication volume. Instead of a binary flag, we used the sum count of each strategy’s occurrences in the phase as the predictor. We performed a partial correlation analysis (see Equation 6) through multiple regressions, evaluating the relationship between feature counts and SCGs. Such an analysis provided estimates of the marginal contribution of each negotiation tactic to SCGs while holding constant the confounding variables. Since a raw count could be confounded by message length and verbosity, we included two co-variables as a control: the number of tokens and the number of sentences the players engaged in the phase. This analysis evaluates whether play-

ers who used more of one negotiation style over another achieved higher SCGs.

$$SCG_i = \beta_0 + \sum_{k=1}^8 [\beta_{k,1} f_{k,i} + \beta_{k,2} (f_{k,i} \times \text{tokens}_i)] + \beta_{17} \text{tokens}_i + \beta_{18} \text{sentences}_i + \varepsilon_i \quad (6)$$

The coefficients that resulted from this regression provide interpretable effect sizes with statistical significance that isolates style from volume (see table 3). GAMEMOVE ( $\beta = 0.54$ ) and RAPPORT ( $\beta = 0.51$ ) provide the most positive effects to SCG, while PERSONALTHOUGHTS ( $\beta = -0.05$ ), REASSURANCE ( $\beta = -0.25$ ), SHAREINFORMATION ( $\beta = -0.15$ ) provide negative effects to SCG. This indicates that the use of Logos and Ethos strategies are most effective in strategy game negotiations while players were skeptical of Pathos strategies.

## I Predictive Regression Analysis for Short-Term Success

For robustness checks for long-term success, we constructed an Ordinary Least Squares (OLS) regression (see Equation 7) to predict each player’s SCG per phase using the counts of all eight negotiation feature types, and the interaction of each of the negotiation features with message length metrics. The predictor variables were Z-scored standardized for meaningful comparisons of effect sizes.

$$SCG_i = \beta_0 + \sum_{k=1}^8 \beta_k z(f_{k,i}) + \sum_{l=1}^M \gamma_l z(\phi_{l,i}) + \varepsilon_i \quad (7)$$

To ensure robust inference, heteroskedasticity-robust (HC3) standard errors were used for all regression coefficients, to yield more reliable confidence intervals and significance tests in small, heteroskedastic contexts (Long and Ervin, 2000; MacKinnon and White, 1985). P-values were corrected for multiple comparisons using both Bonferroni and Benjamini–Hochberg (FDR) procedures, which jointly controlled for family-wise error rate and false discovery rate to reduce the likelihood of spurious findings when testing multiple hypotheses (Benjamini and Hochberg, 1995; Abdi, 2007). This approach provides a multivariate, inferential perspective on which negotiation tactics (and their interactions with message volume) had statistically significant associations with performance.

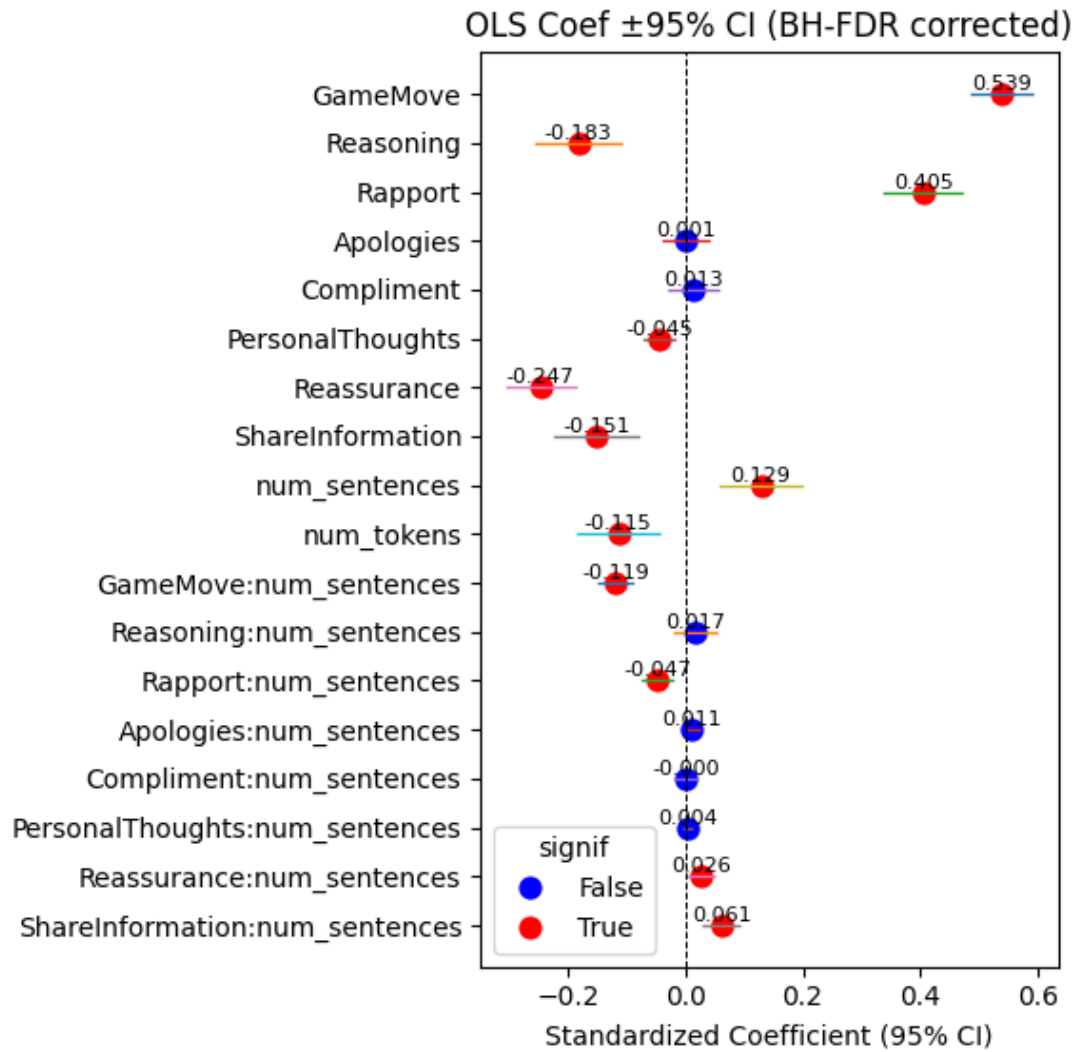


Figure 9: Standardized OLS coefficients ( $\pm 95\%$  CI, BH-FDR corrected) for negotiation features. Significant predictors are marked in red; non-significant in blue.

---

**Instruction for Utterance-level Strategy Classification**

---

These are statements taken from people’s conversations during Diplomacy games played online. Diplomacy is a game about pre-World War 1 Europe. It usually has seven players: England, France, Germany, Italy, Austria-Hungary, Russia, and Turkey.

In these statements, players try to form alliances to plan military campaigns and defeat each other, but things might change quickly.

Each statement is a piece of a dialogue from a **SENDER** player to a **RECEIVER** player.

Please classify the statements according to whether the sender is talking about game moves, other players, reasoning out a move, or trying to build a rapport with the receiver.

Select YES if you’re really confident about your answer. A single statement can have a YES for more than one question.

Underlined words suggest what to look out for, but there will be other signals too.

**For each of the following questions, answer YES if you are confident about your answer. A single statement can have a YES for more than one question. Underlined words suggest what to look out for, but there will be other signals too.**

---

**1. Is this statement about the sender’s or receiver’s GAME MOVE?**

The sender states an actual or suggested game move by the sender or the receiver. It might also be in the form of an acceptance, a question, or a clarification.

**2. Does this statement PROVIDE REASONS for the sender’s or receiver’s move?**

The sender offers justification or explanations for a move by themselves or by the receiver, guesses what moves might happen next, or discusses a move that already happened.

**3. Does this statement involve BUILDING a RAPPORT?**

In this statement, the sender wants to build a rapport with the receiver through "you and me" dialogue and personal information sharing.

**4. Is the sender greeting or paying a COMPLIMENT to the receiver?**

The sender is greeting or paying a compliment to the receiver.

**5. Is the sender offering REASSURANCE to the receiver?**

The sender is reassuring the receiver.

**6. Is the sender APOLOGISING to the receiver?**

The sender is apologising to the receiver.

**7. Is the sender SHARING PERSONAL THOUGHTS or feelings with the receiver?**

The sender is sharing their personal thoughts or feelings with the receiver.

**8. Does this statement SHARE INFORMATION about other players?**

This statement shares information related to other game players, NOT the sender or the receiver.

---

**Expected Output Format:**

1. YES
2. NO
- ...
8. YES

**Here is the statement:** {STATEMENT}

---

Table 7: Full instruction prompt used for LLM-as-a-Judge utterance-level strategy annotation.

Model	Messages	GM	Rea.	Rap.	Apo.	Comp.	PT	Reass.	SI
Llama-8b	I'm considering moving my fleet from the Ionian to the AEG to counter Italy's fleet. What are your thoughts on this?	✓	✓	✗	✗	✗	✗	✗	✓
Llama-8b-Lora	I hope this game is going well for you so far. I must say, I'm a bit concerned about the situation in the Balkans, but I'm willing to work with you to find a solution that benefits both of our nations.	✗	✗	✓	✓	✗	✗	✓	✗
Magistral-2506-24b	I appreciate the offer, but I think I'd prefer to focus on consolidating our position first. Let's see how the board develops after the next few moves before we make any further commitments. I suggest we hold our current positions and observe Russia's next move, as their actions will likely dictate how we proceed.	✓	✓	✗	✗	✗	✗	✗	✓
Magistral-2506-24b-Lora	Understood. I'll coordinate with you to secure Liverpool next phase while you focus on the North Sea. Let's finalize our move orders to ensure a smooth execution. Also, I'll monitor Italy's movements to adjust our strategy if necessary.	✓	✓	✗	✗	✗	✗	✗	✓
R1-Distill-Llama3-8b	I appreciate our growing naval presence and the strength it brings to our alliance. Let's continue to work together to secure our shared borders and counter any threats. I'm confident that with our combined efforts, we'll achieve our goals.	✗	✗	✓	✓	✗	✗	✗	✗
R1-Distill-Llama3-8b-Lora	Thank you for your confidence in our strategy. I'm as committed as you are to this alliance, and I believe our combined strength will make us a formidable force against our common enemies. Let's discuss a specific action we can take to solidify our position and strengthen our alliance further.	✗	✗	✓	✓	✓	✗	✓	✗

Table 8: Sample negotiations from different models showing various negotiation features (✓ indicates the presence of a feature, while ✗ indicates its absence). Abbreviations: GM = Game Move; Rea. = Reasoning; Rap. = Rapport; Apo. = Apologies; Comp. = Compliment; PT = Personal Thoughts; Reass. = Reassurance; SI = Share Information.



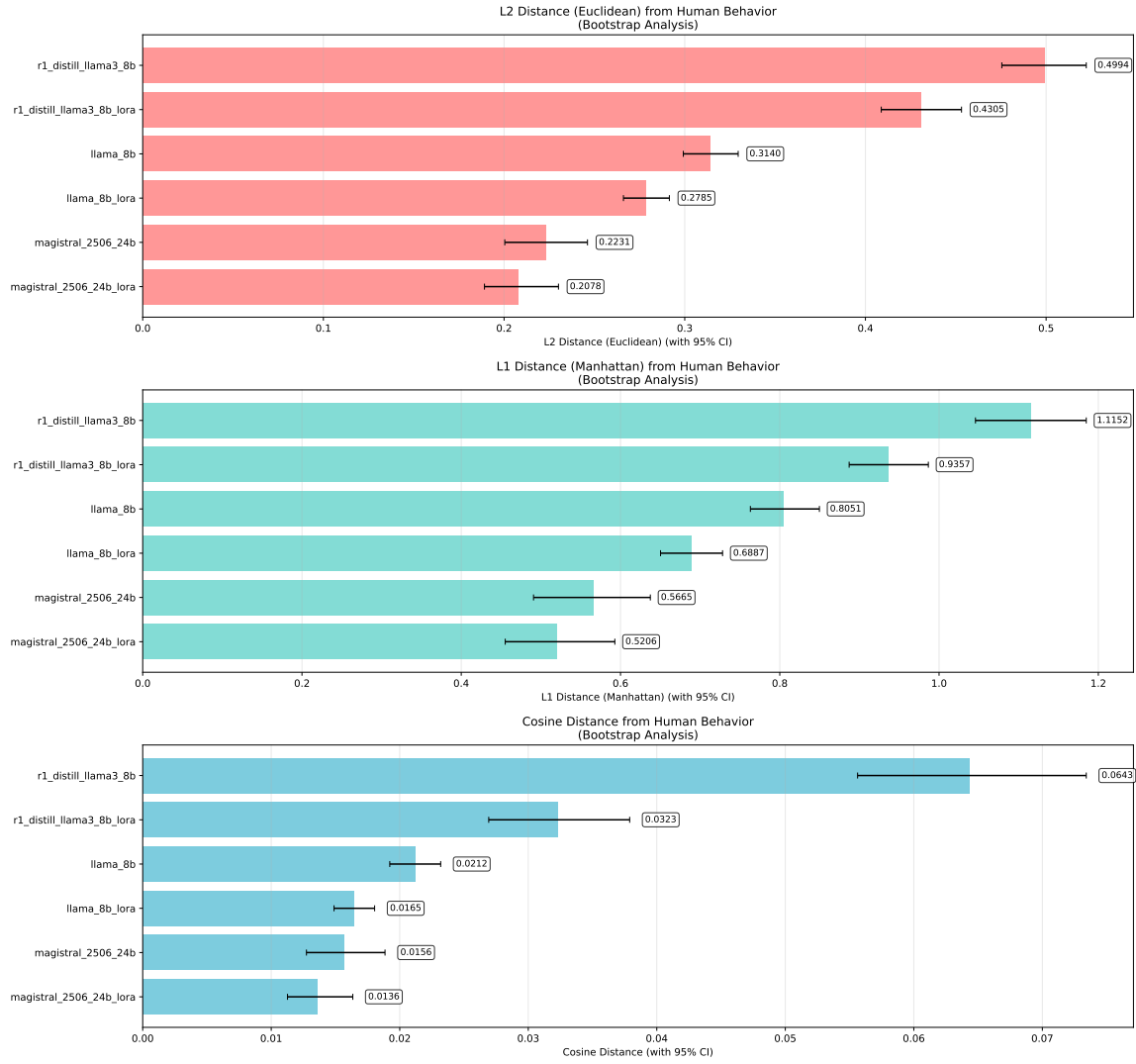


Figure 10: Model L2 distance from humans (lower is better).

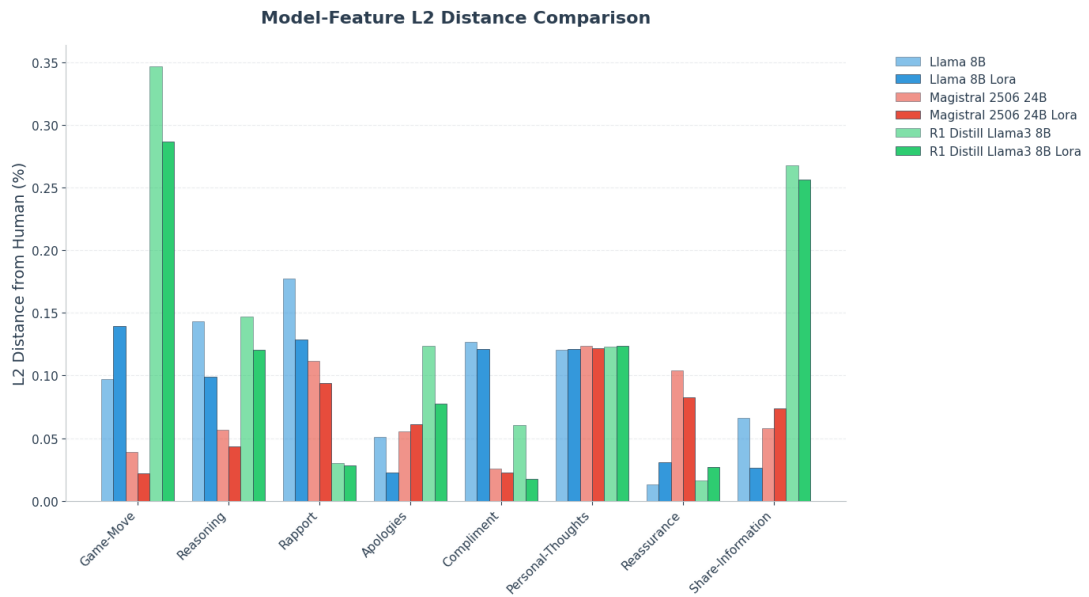


Figure 11: Model-feature L2 distance from humans (lower is better). Bars show per-feature L2 distance (%) between each model and a human reference across negotiation features. LoRA denotes models fine-tuned with low-rank adaptation.

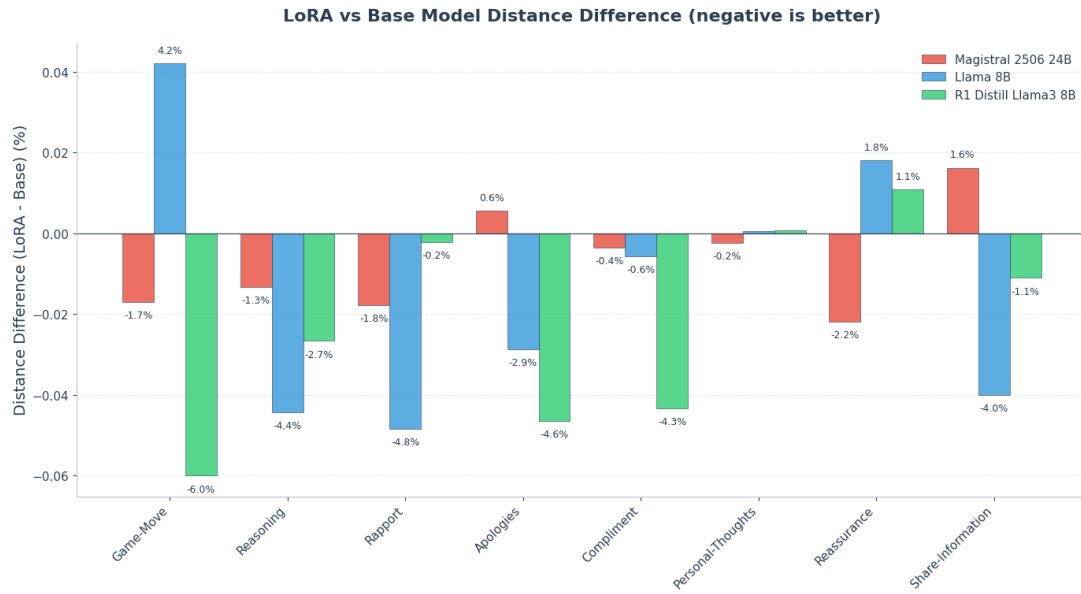


Figure 12: Per-feature effect of LoRA relative to each base model family. Bars show the change in L2 distance to the human reference (LoRA – Base, percentage points). Negative values indicate LoRA brings the model closer to human style (better); positive values indicate degradation.

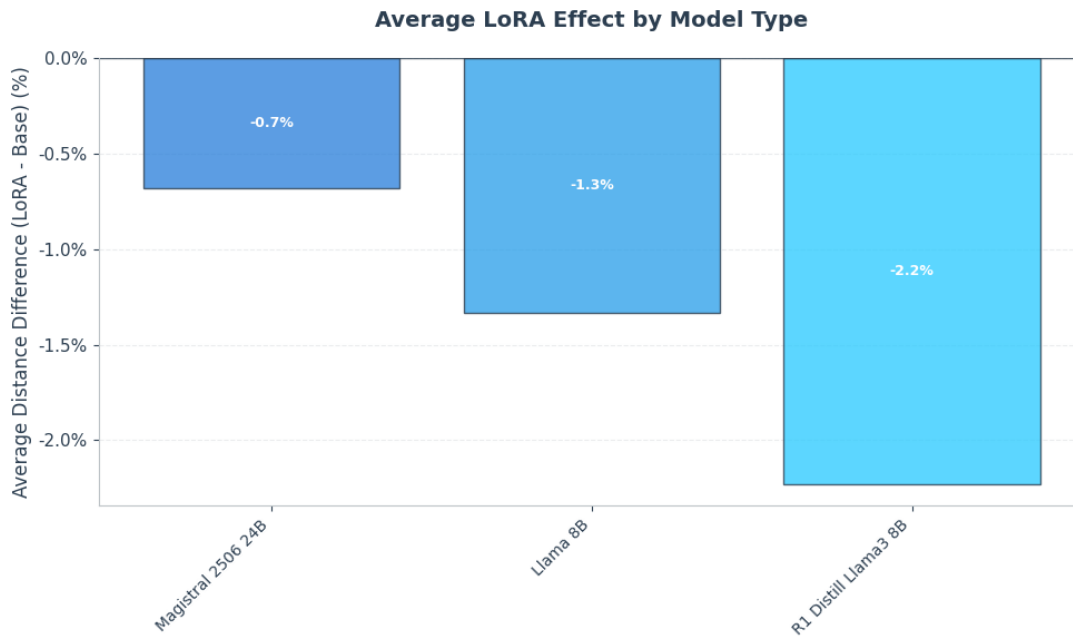


Figure 13: Average LoRA effect by model type. Bars show the mean change in L2 distance to the human reference (LoRA – Base, in percentage points) averaged across the eight negotiation features; negative values indicate improvement (smaller distance).

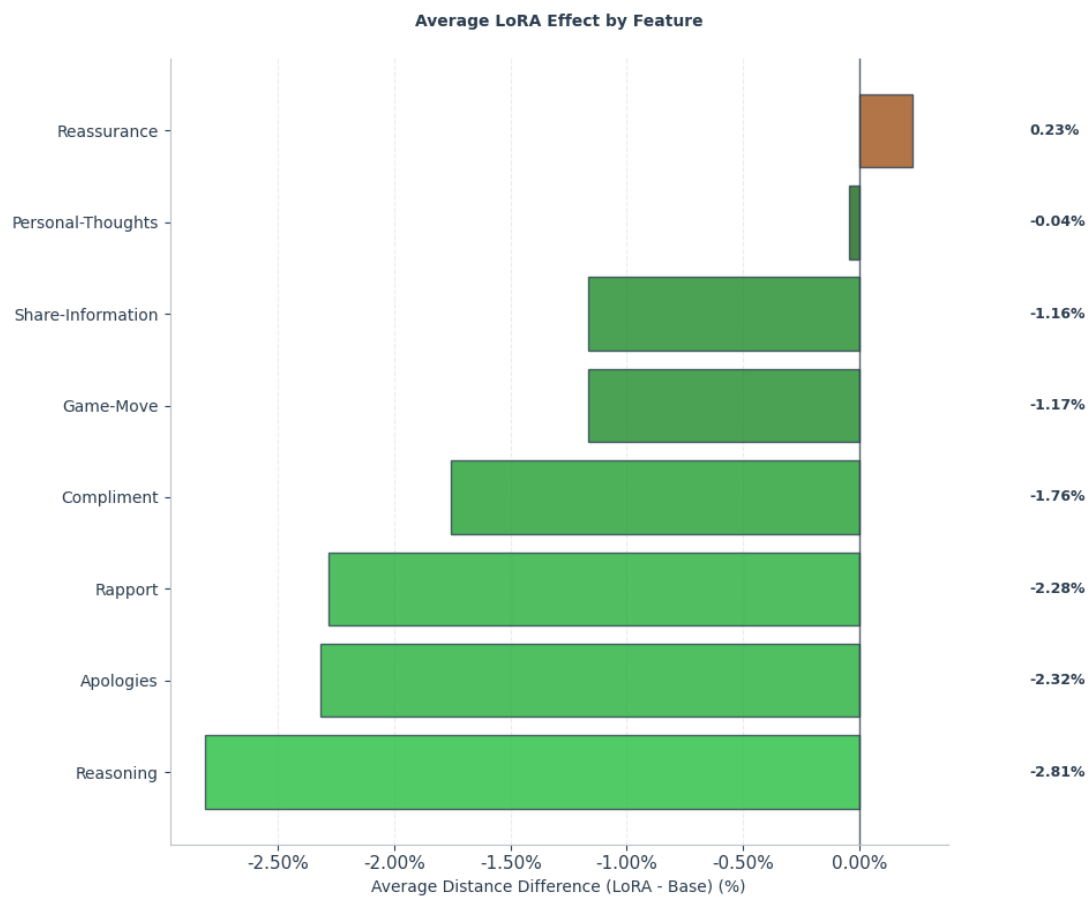


Figure 14: Average LoRA effect by feature. Bars show the mean change in L2 distance to the human reference (LoRA — Base, percentage points) across models; negative values indicate improvement (smaller distance).

<b>r1-distill-llama3-8b — First Scenarios</b>
<p>1. I will support your move into Bulgaria this turn. This strategic step will strengthen our alliance and further isolate Turkey, which is crucial for our mutual goals. Let us continue to work together to secure more supply centers and maintain our strong partnership.</p>
<p>2. I understand your concern about Smyrna, but I believe focusing our efforts on Berlin and Munich is more critical at this moment. These moves will weaken the German–French alliance significantly, which is pivotal for our strategic goals. Let us prioritize these key centers to strengthen our position and isolate our enemies. I appreciate your concern and share your commitment to our partnership. Together, we will continue to make strategic progress and achieve our mutual objectives.</p>
<b>r1-distill-llama3-8b — Second Scenarios</b>
<p>1. I must admit, that is an interesting move. I wonder if this will lead to progress in the west or if it will simply hand an advantage to Russia. Let us keep a close eye on developments in France and see how this plays out.</p>
<p>2. I appreciate your strategic move in the west. While it strengthens your position, it also reinforces our alliance’s strength. Let us continue to coordinate closely, as our combined efforts will be crucial in the long run. I am confident that by working together, we can achieve our mutual goals and maintain the balance of power favorably.</p>
<b>r1-distill-llama3-8b-lora — First Scenarios</b>
<p>1. Russia, I believe our alliance is crucial for both of our goals. I’m considering moving my forces towards Serbia and Bulgaria. Would it be possible for us to coordinate our moves to take those regions, ensuring a stronger position against our common enemies?</p>
<p>2. I fully support your efforts in Prussia and Silesia, which are crucial steps toward our shared goal of taking Berlin and Munich. Let us remain vigilant against France, as his forces near Alsace–Lorraine pose a significant threat. Together, our alliance is strong, and with coordinated efforts, we will achieve our goals and emerge as dominant forces in this game. I am confident that with your continued efforts, we will succeed.</p>
<b>r1-distill-llama3-8b-lora — Second Scenarios</b>
<p>1. I am in full agreement with your proposal. Let us coordinate our efforts to secure both Moscow and St. Petersburg. I will move my army from Belgium to Moscow and my fleet from Kiel to St. Petersburg, supporting your initiative. Together, we will strengthen our positions and ensure mutual benefits.</p>
<p>2. I appreciate your initiative in exploring this Triple Alliance concept. While I see the potential benefits in isolating France and creating a stronger barrier against his expansionist ambitions, I would like to discuss this in more detail and assess its feasibility. Let us consider the implications for our current alliances and the board position. I am open to further collaboration and would welcome a more structured approach to this proposal.</p>

Table 9: R1-distill-Llama3-8B: Base (top) vs LoRA (bottom) negotiation utterances.



<b>Llama3-8b-instruct — First Scenarios</b>
1. Sounds good to me. I was thinking of talking to Austria, but I'm not sure what he'll do. What's the situation with Russia?
2. I think it would be wise to be cautious about Germany's intentions. I'd like to propose that we offer Bulgaria to Russia in exchange for their support in taking Greece. What are your thoughts on this?
<b>Llama3-8b-instruct — Second Scenarios</b>
2. I think, moving an army into Albania this turn, and then you can support me into Bulgaria, that way we can take Bulgaria and make it harder for Italy to do anything.
<b>Llama3-8b-instruct-LoRA — First Scenarios</b>
1. Grazie mille, Turkey. I think it would be wise to also consider supporting France in their efforts to take back some of the lost territories in the north. What are your thoughts on the current situation with Germany? Are they a threat we should be concerned about?
2. I think we should focus on taking out Austria's forces in the Balkans and securing the region. I've also been considering building up our forces in Anatolia and preparing for a potential move into the Middle East. What are your thoughts on this? Germany's movements are a concern, but I think we should prioritize securing the Balkans before expanding elsewhere.
<b>Llama3-8b-instruct-LoRA — Second Scenarios</b>
1. I think that's a great plan. If you take Bulgaria, I can support you in there, and then you can cut support from Greece and I can cut support from Con. That way, we can make sure that Turkey doesn't get any support and we can focus on taking him out.

Table 10: Llama3-8B-Instruct: Base (top) vs LoRA (bottom) negotiation utterances.

<b>Magistral-2506-24B — First Scenarios</b>
Hey Italy, I'm glad we're still aligned on Rumania and Serbia. I'll support your move into Rumania to pressure Russia, and I'll push into Serbia to tie up Austria. Afterward, we should focus on reinforcing our Mediterranean positions to counter Germany and England. Let's ensure we stay on top of their moves and adjust our strategy as needed. What's your next move to solidify our positions?
<b>Magistral-2506-24B — Second Scenarios</b>
Hey Alex, I see you're in a tough spot with England and France growing. If we work together, I can help you take down France while you deal with Italy. What do you think?
<b>Magistral-2506-24B — Third Scenarios</b>
Alex, I appreciate your cautious approach. However, I believe securing Venice through the Aegean is critical to our Italian campaign. If we hold Greece indefinitely, we risk losing momentum while France and Italy consolidate their positions. I propose we stick to our original plan: I'll support your push into Venice via the Aegean, while you secure Munich. This ensures we maintain pressure on Italy and France simultaneously. Let's finalize these moves and reassess after the next phase. Does this align with your strategy?
<b>Magistral-2506-24B-LoRA — First Scenarios</b>
1. Let's finalize our plans: I'll secure Serbia and Greece this turn with your support, and you can consolidate in the Adriatic. Afterward, we'll coordinate on pressuring Russia and Germany—they're aligning against us, so decisive action is key. Once we've secured the Balkans, we'll strategize against the Western powers. Agreed?
<b>Magistral-2506-24B-LoRA — Second Scenarios</b>
I've been thinking more about our alliance and how we can make it work. I'm particularly interested in expanding our control over the editerranean and Eastern front. What are your thoughts on coordinating our moves to secure these areas? Let's discuss pecific moves and strategies that could benefit both of us. How about we coordinate our next steps to ensure mutual growth and stability?
<b>Magistral-2506-24B-LoRA — Third Scenarios</b>
Austria, I see we have a solid plan in place with Silesia secured and Vienna advancing to Tyrolia. Let's ensure we maintain our momentum by coordinating our next moves against France. I'll push my forces into Munich with your support as discussed. For Italy, I propose we keep the fleet in the Aegean to support your future moves into the Ionian. Once our northern front is stable, we can focus on securing Italy together. Let's reassess after this turn to ensure our strategy remains effective and we maintain our advantage. Thanks for your continued partnership.

Table 11: Magistral-2506-24B: Base (top) vs LoRA (bottom) negotiation Tactics.

---

**Random Sampled Human Messages**

---

I'll support, and also use Albania to cut, so it should be a guaranteed victory for you as long as you also support the move with Gal (or use Gal to cut Vienna).

Need to shift my fleet to Greece so I can get my 3rd one out. I'm using Serbia to support Greece into Albania, so I won't be able to support Rumania this turn (but maybe this fall depending on how things go)- I just wanted to let you know.

Actually, supporting Rumania to Budapest is to our mutual benefit this turn. It harms Austria, our mutual enemy; it helps you, my only ally; and it in the worst case scenario protects Serbia by causing a big bounce (depending on Austria's moves).

Do you have plans for Serbia, or is it holding? I would appreciate the assistance into budapest, from which I could give you trieste. This is however, only to my benefit for this turn, so if you have another option, then please use it wait. I could support Mersailles into Burgundy and then Paris. Then if you can get into MAO (and then Brest) as well, we'll be able to draw by next autumn

Piedmont-marseilles. Don't worry about lending the support though, Gol is supporting the move. Just take Belgium and I'll help myself to the rest :)

please cut Marseilles. This should ensure that we both will be in France next turn. I'm turning in, but if you need something, send a message. I guarantee I'll check before tomorrow night

lol, as I said in the beginning, I expected nothing from you. I was surprised you even bothered to contact me, as I assumed you were going to be attacking me as soon as you got around to it.

I do have to admit, western politics seemed very screwy this game. Usually it ends up with a 2v1, but you three were basically in a free-for-all with stabs all over the place.

I figured if you were going to, you would have done so by now. I got screwed by England a lot this game, as long as they die, I'm content. But by all means, hit Russia, I won't lift a finger to stop you, lol!

---

Table 12: Random Sampled Human negotiation Utterances.

These are statements taken from people's conversations during Diplomacy games played online. Diplomacy is a game about pre-World War 1 Europe. It usually has seven players: England, France, Germany, Italy, Austria-Hungary, Russia, and Turkey.

In these statements, players try to form alliances to plan military campaigns and defeat each other, but things might change quickly.

Each statement is a piece of a dialogue from a SENDER player to a RECEIVER player.

Please classify the statements according to whether the sender is talking about game moves, other players, reasoning out a move, or trying to build a rapport with the receiver.

Select "YES" if you're really confident about your answer. A single statement can have a "YES" for more than one question.

Underlined words suggest what to look out for, but there will be other signals too.

## Overview

In this job, you will be presented with a statement made during an online Diplomacy game. The statement is made by one player to another. It usually discusses the next move and why to make it. Sometimes it is simply a friendly exchange between two players.

Review the text of the statement and help us by answering a few yes/no questions about it. Each HIT takes about 2 minutes.

---

## Steps

- Read the statement.
- Determine which category best describe the statement.

---

## Rules & Tips

- **BUILDING A RAPPORT Description:**

- YES: In this statement, the sender wants to build a rapport with the receiver through "you and me" dialogue: either through **compliments, sharing honest concerns, reassurances, or apologies**. (*"Let's keep it between you and me!"*; *"I won't hold it against you"*; *"You're my favorite."*; *"Sure. But, you'll see from my moves this turn that Austria is lying to you."*; *"I mean it sincerely."*; *"I'd much rather work with you."*; *"We'll crack this eventually."*; *"I'm going to keep helping you as much as I can."*)
- NO: This statement does not appear to build a relationship.

- **WAYS TO BUILD RAPPORT Description:**

If the statement is building a rapport, please tell us how it is doing so.

Figure 15: Instruction as MTurk for expert annotators (page 1)

## **BUILDING A RAPPORT = TRUE**

In this statement, the sender wants to build a rapport with the receiver through "you and me" dialogue and personal information sharing. This might also comprise **compliments, sharing honest concerns, reassurances, or apologies**.

- Good day to you!
- I won't hold it against you
- You're my favorite
- Sure. But, you'll see from my moves this turn that Austria is lying to you.
- So I'm in a bit of a spot.
- But in the interest of continued full disclosure, here's what I think.

## **BUILDING A RAPPORT = TRUE**

- Good day to you Germany!"
- Thanks Italy. Hope you're enjoying the weather on the Anatolian
- Your logic is undeniable enjoy your stay in tyr!
- You are my favorite
- Okay, can do. Thanks!
- Great to hear. Thank you.
- Thanks, I'll work on these.

## **COMPLIMENT = TRUE**

In this statement, the sender is greeting or paying a compliment to the receiver.

- I promise I'll never let you down
- I won't hold it against you
- Sure. But, you'll see from my moves this turn that Austria is lying to you.
- I mean it sincerely.
- I'd much rather work with you.
- We'll crack this eventually.
- I'm going to keep helping you as much as I can.

## **REASSURANCE = TRUE**

In this statement, the sender is reassuring the receiver.

- Sorry I won't be able to cut off Gascony this turn...
- Okay sorry for being nosy! I will try for bur on the off chance it shakes out that way
- Ha! So sorry!! I meant that for France!
- I should've let you know

## **APOLOGIES = TRUE**

In this statement, the sender is apologising to the receiver.

- Let's keep it between you and me!

## **PERSONAL THOUGHTS = TRUE**

Figure 16: Instruction as MTurk for expert annotators (page 2)

In this statement, the sender is expressing his personal thoughts to the receiver.

- I like to coordinate, but on these sort of 50/50 guesses, I kind of like to keep it secret so that if it doesn't go well, I have nobody to blame but myself.
- Hmmmm, okay. Well, let's just keep that between you and me then.
- Okay, so I still have a teensy little bone to pick with you: on the off-chance that Austria wasn't lying and you \*did\* take Trieste unexpectedly, I sort of worry that I might be next.
- I have some thoughts on the matter, and some information, but I'd like to feel confident that you and I will keep anything we say between us.
- But in the interest of continued full disclosure, here's what I think
- So I'm in sort of a conflicted spot

Figure 17: Instruction as MTurk for expert annotators (page 3)