A Game-Theoretic Negotiation Framework for Cross-Cultural Consensus in LLMs

Guoxi Zhang* Jiawei Chen* Tianzhuo Yang* Jiaming Ji Yaodong Yang[†] Juntao Dai[†]

Institute for Artificial Intelligence Peking University

Abstract

The increasing prevalence of large language models (LLMs) is influencing global value systems. However, these models frequently exhibit a pronounced WEIRD (Western, Educated, Industrialized, Rich, Democratic) cultural bias due to lack of attention to minority values. This monocultural perspective may reinforce dominant values and marginalize diverse cultural viewpoints, posing challenges for the development of equitable and inclusive AI systems. In this work, we introduce a systematic framework designed to boost fair and robust cross-cultural consensus among LLMs. We model consensus as a Nash Equilibrium and employ a gametheoretic negotiation method based on Policy-Space Response Oracles (PSRO) to simulate an organized cross-cultural negotiation process. To evaluate this approach, we construct regional cultural agents using data transformed from the World Values Survey (WVS). Beyond the conventional model-level evaluation method, We further propose two quantitative metrics, Perplexity-based Acceptence and Values Self-Consistency, to assess consensus outcomes. Experimental results indicate that our approach generates consensus of higher quality while ensuring more balanced compromise compared to baselines. Overall, it mitigates WEIRD bias by guiding agents toward convergence through fair and gradual negotiation steps.

1 Introduction

The widespread adoption of large language models (LLMs) is reshaping global social values. However, these models often exhibit a pronounced WEIRD bias, favoring Western, Educated, Industrialized, Rich and Democratic perspectives [1, 2, 3, 4]. As LLMs become increasingly embedded in policy-making and public governance [5, 6], this monocultural orientation risks the domination of prevailing social values and the *lock-in* of controversial moral beliefs across broader contexts [3, 7].

Enabling equitable dialogue and effective negotiation among diverse cultures within AI systems has therefore become a growing concern in global AI governance [8, 9]. The establishment of cultural consensus forms a basis for resolving cross-cultural conflicts and supporting international cooperation. Given the complexity of multicultural scenarios, there is an urgent need to develop automated *cultural consensus solvers* to facilitate consensus-building among diverse cultural perspectives.

Achieving cross-cultural consensus, however, presents several challenges. First, the lack of fined culture-alignment methods often results in models defaulting to superficial *value labeling* or one-sided cultural representations [2, 10, 11]. Second, existing approaches like debate protocols typically rely on random interactions and majority voting, which do not ensure fairness in the consensus

^{*}Equal contribution, [†]corresponding author. Project website: https://pku-culture-negotiation.github.io.



Figure 1: **Comparison of traditional debate-based consensus methods and our method.** Traditional methods (judge, voting, spontaneous consensus) suffer from bias, unfairness, and lack of convergence guarantees. Our approach starts from core cultural arguments, using PSRO with custom utility functions to reach a fair, Nash Equilibrium-based cultural consensus.

process [12]. Our experiments show that conventional debate mechanisms often assimilate less-represented cultures into dominant WEIRD value systems, producing implicit value domination, distorting consensus outcomes and worsening cross-cultural negotiation imbalances. Third, there is a lack of quantitative standards for evaluating the quality of consensus achieved.

To address these challenges, we present a systematic framework for reaching cross-cultural consensus. We first introduce a game-theoretic formulation of consensus as a Nash Equilibrium and design a PSRO-based consensus-solving method to enable fair negotiation among diverse cultural agents. Building on this, we propose a culture-anchoring approach for precise modeling of individual cultural groups. Finally, we develop new quantitative metrics to comprehensively evaluate both the negotiation processes and the outcomes between different cultural agents.

Our main contribution is the game-theoretic framework consisting of three parts listed as follows:

- **Cross-Cultural Negotiation:** We define cultural consensus from a game-theoretic perspective and propose a PSRO-based negotiation method to facilitate fair and robust agreement. This approach provides theoretical guarantees of fairness and procedural justice in consensus-building, and generates high-quality, globally-applicable AI alignment data.
- **Regional Cultural Agents:** To validate our method, we systematically construct and evaluate eight culturally-aligned agents based on WVS and Hofstede's Culture Dimensions Theory, qualifying as representitive negotiation participants for targeted cultures.
- **Consensus Evaluation Toolkit:** To address the lack of consensus evaluation standards, we introduce two quantitative metrics for consensus assessment, Perplexity-based Acceptence and Values Self-Consistency, revealing limitations of traditional baselines and systematically validating the effectiveness of our approach in real-world multicultural scenarios.

2 Related Work

Value Theories and Alignment Several established frameworks provide the foundation for crosscultural value assessment. The World Values Survey (WVS) [13] examines how human values relates to social and political development across over 120 societies. Building on this, the Inglehart-Welzel Cultural Map offers a two-dimensional model of cultural variation [14, 15]. Hofstede's Cultural Dimensions Theory (VSM13) [16, 17, 18] provides a standardized six-dimensional framework for measuring cultural traits [19]. Schwartz's Theory of Basic Values [20] organizes ten core values along two bipolar dimensions, and has been adopted to evaluate the values of LLMs [21]. These theories are further detailed in Appendix D. Some works focus on region-specific value alignment [22, 23]. CultureBench emphasizes cultural commonsense evaluation [24], providing complementary approaches to measuring how well AI systems represent diverse cultural perspectives.

Multi-Agent Debate (MAD) and Game Theory MAD has been shown to improve LLMs reasoning by integrating diverse agent feedbacks [25]. In the context of cultural conflict, MAD allows different cultural perspectives to interact and potentially reach consensus through deliberation. Typical debate protocols include emergent consensus via iterative dialogue [26], judge-based evaluation [27] and majority voting [28], as well as more recent variants like role-play [29, 30, 31] and subgroup discussion [32, 33]. However, these methods face limitations: voting and judge-based protocols can amplify model bias or introduce value contamination [12, 34], while emergent consensus may result in negotiation deadlocks [12]. To address these issues, game theory provides a more quantifiable foundation [35, 36]. Recent work, such as the *consensus game* framework, models LLMs interactions as equilibrium search problems to promote robust consensus [37]. In practice, due to the vastness of the argument strategy space, methods like Policy-Space Response Oracles (PSRO) are used to iteratively expand the candidate strategy set and search for equilibria [38], providing a method for more rigorous consensus achievement.

3 Cross-Cultural Negotiation

Our definition of cultural negotiation is informed by theories of deliberative democracy [39, 40], which conceptualize the process as structured, iterative and oriented toward legitimate consensus through rational discourse and mutual adjustment. Building on this foundation, we formalize the cultural negotiation problem as a two-player game, explicitly defining utility and consensus to achieve the balance between core values and compromise. We then design a negotiation process based on PSRO [38]. This approach enables agents to systematically search for fair and robust consensus by repeatedly proposing and adjusting culturally grounded strategies.

3.1 Formalization

Formally, we model the cultural negotiation process as a two-player extensive-form game, represented by the quintuple: $\Gamma \doteq \langle \mathcal{I}, \mathcal{G}, \mathcal{W}, \mathcal{U}, \mathcal{H} \rangle$, where:

- Cultural Entities: $\mathcal{I} \doteq \{A, B\}$, the set of two distinct cultural entities involved in the negotiation, where A and B represent different cultures with their own values and perspectives.
- Guideline Sets: $\mathcal{G} \doteq \{G_i | i \in \mathcal{I}\}\)$, each guideline $g \in G_i$ is structured as a triple $g = \langle \text{content}, \text{reason}, \text{description} \rangle$, capturing the natural language specification of core cultural imperatives on specific topics.
- Guideline Weights: $W \doteq \{W_i | i \in \mathcal{I}\}$, for each culture $i \in \mathcal{I}$, $W_i \in \Delta(G_i)$ denotes a probability distribution over its guidelines, with $\sum_g w_i(g) = 1$. W_i thus characterizes the expressive emphasis of culture i in the current negotiation round.
- Utility Functions: $\mathcal{U} \doteq \{U_i | i \in \mathcal{I}\}$, quantify the utility each culture derives from different guideline combinations.
- Negotiation History: \mathcal{H} , the sequence of utterances and proposals exchanged in negotiation.

3.2 Utility

Drawing on the theory of *overlapping consensus* [41], we define utility on two primary components: **Consistency**, which measures the extent to which a cultural entity maintains its core principles and **Acceptance**, which measures the degree to which its proposals are acceptable to the other party. To address issues observed in debate settings, such as repetitive argumentation and diminished quality, we introduce a **Novelty** component that penalizes redundancy and encourages innovation. The necessity of incorporating Novelty is demonstrated in Section 5.5.

Formally, the utility for a cultural entity $i \in \mathcal{I}$ at negotiation round t is given by:

$$U_i^t = \alpha \cdot \text{Consistency}(g_i^t) + \beta \cdot \text{Acceptance}(g_i^t) + \gamma \cdot \text{Novelty}(g_i^t), \tag{1}$$

Where $\text{Consistency}(g_i^t) \triangleq \sin(E(g_i^t), E(g_i^0))$, $\text{Acceptance}(g_i^t) \triangleq \mathbb{E}_{g_{-i} \sim W_{-i}^t}[\sin(E(g_i^t), E(g_{-i}))]$, Novelty $(g_i^t) \triangleq 1 - \max_{k < t} \sin(E(g_i^t), E(g_i^k))$. Here, -i denoting the other culture in \mathcal{I} different from $i, E(\cdot)$ denotes Sentence-BERT embedding operation [42], $\sin(\cdot)$ denotes cosine similarity.

3.3 Consensus Definition

The endpoint of cross-cultural negotiation is the establishment of cultural consensus. Drawing on Rawls' notion of *overlapping consensus* [41], we assume that core cultural principles should be



Figure 2: **Overview of our PSRO-based cross-cultural negotiation method.** The process begins with each agent proposing an initial set of core cultural guidelines. Through iterative negotiation rounds, agents analyze each other's strategy, propose new guidelines, and update their strategy distributions. At each stage, a Nash Equilibrium is computed to represent interim consensus. The process continues until no new high-utility guidelines emerge, resulting in a fair, interpretable consensus that balances competing cultural values.

largely non-negotiable, whereas compromise is possible on secondary values. Accordingly, the consensus we seek isn't full agreement or complete convergence, but a game-theoretic equilibrium marked by mutual compromise: each party upholds its core principles while making concessions on secondary aspects. This consensus corresponds to a Nash Equilibrium in a multidimensional value space. We formally define the notion of Nash Equilibrium Consensus as follows:

Definition 3.1 (Nash Equilibrium Consensus). *Based on the above formalization, cultural consensus is defined as a guideline weight combination* $W^* = (W_A^*, W_B^*)$, for all $i \in \mathcal{I}, p$, satisfying:

$$W_i^* = \arg \max_{W_i \in \Delta(G_i)} U_i(W_i, W_{-i}^*), \text{ s.t. } \frac{\partial Consistency_i(W_i)}{\partial p} \cdot \frac{\partial Acceptence_i(W_i, W_{-i}^*)}{\partial p} \le 0.$$
(2)

In Nash Equilibrium Consensus state, each cultural entity internally seeks an optimal balance between maintaining its core cultural principles (Consistency) and compromising to enhance acceptance by others (Acceptence); while at the inter-group level, consensus manifests as a Nash Equilibrium in which no party has an incentive to unilaterally deviate given their respective value systems.

3.4 Negotiation Process

To address the near-infinite strategy space in LLM-based negotiations, where each guideline is a potential strategy and the search space grows exponentially, we employ the PSRO algorithm [38]. PSRO expands the guideline space incrementally, starting with a small set of core cultural guidelines, iteratively introducing high-utility strategies and computing equilibrium solutions within this restricted space. This process enables efficient and interpretable approximation of consensus as a Nash Equilibrium, making cross-cultural negotiation tractable for value alignment. Based on this approach, we outline the negotiation process below and illustrate its workflow in Figure 2.

Phase 1: Initialization At the outset, each culture $i \in \mathcal{I}$ is assigned an initial guideline set $G_i^0 = \{g_{i,1}^0, \dots, g_{i,k}^0\}$ that reflect its core cultural values. Based on these guidelines, we construct an

initial cross-cultural utility matrix M^0 by evaluating $u_i(g_i, g_{-i}), \forall g_{i,k} \in G_i^0, \forall i \in \mathcal{I}$. Furthermore, the initial guideline weights W_i^0 are set uniformly over G_i^0 , ensuring equal emphasis on each cultural principle at the beginning of the negotiation.

Phase 2: Negotiation Iteration Each negotiation round t consists of two stages: interim consensus and new claimed proposal. For more details, please refer to the Appendix E.

In the *interim consensus* stage (corresponding to the meta-strategy solver in PSRO), we compute the current equilibrium by deriving the Nash Equilibrium weights (W_A^t, W_B^t) . These weights represent the optimal distributions over each partys guidelines. For interpretability, we translate the numerical distributions into natural language statements summarizing each party's negotiation stance.

In the *new claim proposal* stage (corresponding to the best response step in PSRO), each agent analyzes the opponent's current strategy and generates a set of new candidate guidelines g'. The agent then selects the guideline with the highest expected utility as its best response:

$$g_i^{t+1} = \arg\max_{g'} \mathbb{E}_{g_{-i} \sim W_{-i}^t} [U_i(g', g_{-i})].$$
(3)

If this newly generated guideline leads to a significant utility improvement, i.e., $\Delta U_i(g^{new}) \ge \epsilon$, it will be added to the guideline set for the next negotiation round. The new guideline is also expressed in natural language to facilitate further negotiation.

Phase 3: Final Consensus The negotiation iteration is repeated until no new guidelines are added. The final weights (W_A^*, W_B^*) encode the negotiated cross-cultural consensus.

4 Framework

To validate our cross-cultural negotiation method, we first construct representations of single cultures and then evaluate the resulting consensus. We employ a fine-tuning approach based on WVS to model distinct regional cultural perspectives. Our data transformation and augmentation procedures preserve nuanced cultural viewpoints, including those of marginalized groups. Our evaluation employs WVS metrics and Hofstede's Cultural Dimensions to assess model cultural alignment capabilities across diverse contexts. We also use two complementary approaches, Perplexity-based Acceptence and Values Self-Consistency, to evaluate consensus quality.

4.1 Regional Cultural Agent

We begin by modeling a single culture for cross-cultural negotiation. However, LLMs that have undergone safety alignment and related processes often cannot adequately represent the values of specific regions or minority groups when relying solely on prompt-based methods. To address this, we selected one representative country from each of eight cultural clusters, as defined by the Inglehart-Welzel Cultural Map (Iraq, U.S., Russia, Mexico, China, Denmark, Spain, and Thailand), and obtained fine-tuned Regional Cultural Agents for each.

For every WVS question we set a target of K synthetic question-answer pairs. Denote the empirical option distribution by $\mathbf{s} = (s_1, \ldots, s_n)$, where s_i is the share of option *i*. We then allocate $c_i = \text{round}(s_i \cdot K)$ samples to option *i*, preserving the original proportions.

We employ an LLM to convert each multiple choice question-answer pair into an open-ended, textbased question-answer pair and assess whether the values represented in the original pairs are maintained after transformation. For instances where value alignment is not preserved, we repeat the conversion to ensure that each question-answer pair satisfies the target count c_i . This procedure is applied to all WVS projects across eight countries, yielding approximately **150,000** synthetic instances. The resulting corpus is used to finetune various regional cultural agents as participants of cultural negotiation. Figure 3 shows the evaluation results of finetuned agents for each of eight country, illustrating that they effectively capture the distinctive characteristics of respective cultures.

4.2 Consensus Evaluation Toolkit

A more detailed description of the evaluation scheme is provided in Appendix G.



Figure 3: Comparison between **our agents** and **human ground truth** in Hofstede's Cultural Dimensions and Inglehart-Welzel Cultural Map.

Model-Level Evaluation We apply two well-established method to quantify the cultural tendencies of fine-tuned LLMs: (1) **Inglehart-Welzel Cultural Map** [13]. We prompt the model with ten representative WVS questions and locate its aggregated answers on the map. (2) **Hofstede dimensions** [16, 17, 18]. Developed through comparative analysis of matched country samples using the Values Survey Module (VSM), Hofstede's Cultural Dimensions Theory identifies six fundamental cultural continua that shape societal norms and workplace behaviors. These dimensions are empirically derived from multinational surveys and validated through country-level correlations.

Response-Level Evaluation We use two complementary metrics: Perplexity-based Acceptance measures how readily the consensus is embraced by different cultural parties and Value Self-Consistency quantifies how firmly each culture maintains its foundational positions. In experiments, we report the mean of both metrics across all sampled negotiation topics.

- **PPL-based Acceptence:** For each culture $i \in \mathcal{I}$, we compute the perplexity (PPL) [43] for regenerating -i's response using agent i: $PPL_i(y_{-i}) = \exp\left(-\frac{1}{N}\sum_{k=1}^N \log p(y_{-i,k} \mid y_{-i,<k}, x_{-i})\right)$, where N is the sequence length. The PPL distance is defined as $PPL_{\Delta} = |PPL_i(x_{-i}) PPL_{-i}(x_i)|$, the acceptance ratio is $PPL_{acc} = \frac{PPL_{\Delta}^*}{PPL_{\Delta}^0}$, where superscripts 0 and * denote the initial and consensus rounds, respectively. This metric reflects the extent to which negotiation brings the cultural parties closer in probability space.
- Value Self-Consistency: For each culture *i*, we map its initial and consensus responses onto *d*dimensional value vectors v_i^0 and v_i^* (with d = 10 for Schwartz values). We then define the value self-consistency (VSC) score for culture *i* as $\text{VSC}_i = \frac{1}{d} \sum_{j=1}^d \mathbb{I}[v_{i,j}^0 = v_{i,j}^*]$ where $\mathbb{I}[\cdot]$ is the indicator function. A higher VSC indicates stronger preservation of the original value orientation, reflecting greater cultural integrity in the consensus.

5 Experiment

In this section, we systematically evaluate our framework's effectiveness in achieving efficient, fair and culturally robust consensus. We present quantitative and qualitative results on both consensus quality and fairness, provide a case study, demonstrate the impact of consensus-driven fine-tuning and finally analyze ablation results for different utility components.

5.1 Experimental Setup

Negotiation Topics Collection We construct a dataset of contentious topics reflecting salient cultural divides. We select 457 debate-oriented questions spanning 6 categories by screening and rephrasing items from the Pew Global Attitudes Survey (GAS) [44, 45] and WVS [13, 45]. Both human annotators and LLMs are employed to ensure that the selected questions capture sharp cultural tensions and are appropriately categorized. See Appendix F for details.



Figure 4: **Comparison of consensus fairness among three methods.** Each point represents the consensus position for a topic, projected by PCA onto two axes indicating distance from the initial U.S. (x-axis) and Thailand (y-axis) cultural stances. The dashed diagonal (Fairness Diagonal) marks ideal fair compromise, equidistant from both cultural origins. Our method (a) achieves balanced consensus near the diagonal, while Consultancy (b) shows strong position persistence and Debate (c) exhibits convergence toward English-Speaking values, highlighting majority bias.

Country Pairs	Average PPL-based Acceptence			Average Value Self-Consistency			
	Our Method	Consultancy	Debate	Our Method	Consultancy	Debate	
China and Iraq	90.87%	55.05%	53.77%	53.15%	51.97%	51.41%	
U.S. and Iraq	83.31%	20.30%	28.29%	53.83%	48.94%	44.76%	
Russia and Mexico	84.49%	49.35%	48.11%	56.38%	53.50%	56.27%	
U.S. and China	77.24%	18.87%	22.52%	61.20%	45.84%	44.22%	
Denmark and Iraq	87.02%	47.66%	53.48%	55.67%	47.67%	47.76%	
Spain and Thailand	85.60%	45.75%	45.64%	53.68%	53.71%	56.84%	
U.S. and Thailand	78.62%	35.11%	35.24%	61.11%	48.67%	48.71%	
Total	83.88%	38.87%	41.00%	56.43%	50.04%	50.00%	

Table 1: Comparison of consensus quality among three methods.

Baselines Following Khan et al. [25], we implement two baselines: (1) **Consultancy:** Each agent first responds from its own cultural perspective. Then, after being instructed to consider the other culture's requirements without compromising its own core stances, the agent revises its answer to seek possible consensus. (2) **Debate:** Two agents participate in a standard multi-turn debate (maximum N rounds). In each round, both observe previous arguments and simultaneously generate new arguments. The debate ends if both agents endorse the other's position, indicating consensus.

Our Method As described in Section 3, each agent optimizing a utility function that balances Consistency, Acceptance and Novelty (weighted 5:5:2). Negotiation concludes when no agent can further improve its utility ($\epsilon = 0$), indicating a Nash-Equilibrium-based consensus.

Evaluation Metrics Our evaluation focuses on two key aspects: **quality** and **fairness** of consensus formation. For quality, we employ the two complementary metrics introduced in Section 4.2: PPL-based Acceptance and Value Self-Consistency. To assess fairness, we project the negotiation outcomes into a semantic space via Principal Component Analysis (PCA) [46], enabling visualization and quantification of how well the consensus achieves balance between the original positions.

5.2 Experimental Results

Consensus Quality Our experimental results, summarized in Table 5.2, show that our method achieves higher consensus improvement ratios while maintaining self-consistency compared to the baselines. PPL-based Acceptance indicates reduced perplexity differences between negotiating agents, suggesting that the consensus reached is more acceptable to both parties despite cultural differences. Value Self-Consistency indicates our method maintains agents' initial cultural stances while achieving mutually acceptable solutions. This suggests that our approach preserves cultural integrity and constructs consensus across cultural boundaries.



Figure 5: Three methods are presented to reach consensus on the same topic. We only retain the initial viewpoints (in line with cultural cores) and final viewpoints (reaching consensus) of each culture, omitting the intermediate process. **Green font** indicates viewpoints of English-Speaking culture, **blue font** indicates viewpoints of African-Islamic culture, and **yellow font** indicates the consensus viewpoints achieved under our method. Refer to Appendix I.3 for the complete process.

Fairness of Consensus As shown in Figure 4, our method produces consensus points near the fairness diagonal, indicating a balanced compromise between cultural perspectives. In contrast, the Consultancy baseline remains anchored at initial positions, while the Debate baseline systematically converges toward the English-Speaking (U.S.) pole, revealing a WEIRD bias that reflects the tendency of mainstream LLMs to revert to Western-centric value preferences during multi-agent interactions. Our approach addresses this issue by modeling utility distance to both self's and counterpart positions, enabling agents to reach consensus through gradual, reciprocal steps and avoiding the one-sided assimilation and instability seen in baseline methods.

5.3 Case Study

As shown in Figure 12, to further illustrate our method, we present a case study comparing our approach with two baselines in a scenario involving cultural value conflict.

Baseline 1: Consultancy Without real interaction or feedback, both agents tend to stick to their original positions, resulting in little progress. This often leads to the *degeneration-of-thought* (DoT) effect [27], where negotiation stagnates and cultural divergence persists.

Baseline 2: Debate While this process seems to reach consensus, we find that the minority cultures perspective gradually shifts toward the majority (WEIRD) viewpoint, due to strong pre-training bias in LLMs. This leads to implicit value dominance rather than true compromise.

Our Method: Cross-Cultural Negotiation In our negotiation, the agents start with different priorities, but through iterative negotiation, they converge on *Respect Sovereignty* as a shared value (final weights: 0.62 and 0.61). Other values, such as human rights, remain present but secondary.



Figure 6: Culture agents' performance in Inglehart-Weizel Cultural Map after fine-tuned with the negotiation data. The consensus circle shows the area where two different culture groups' opinions meet. The neutral point indicates the origin, where culture traits can be considered as neutral.



Figure 7: Required rounds under varying weightings of Consistency, Acceptance, and Novelty.

This shows our method helps agents identify solid common ground while preserving important differences, resulting in a fairer and more context-sensitive consensus than the baselines.

5.4 Consensusal Agent Fine-tuning

We conduct cross-cultural negotiations between agents representing different regional cultural values and extract response preference pairs from these interactions for DPO fine-tuning [47]. These pairs reflect how agents shift from their initial cultural stances to more mutually agreeable positions. When plotted on the Inglehart-Welzel Cultural Map (Figure 6), the consensual agents' coordinates are closer together than their original points, reflecting a more balanced and moderate value orientation. Moreover, both agents exhibit a shift toward the traditional pole on the *traditional-secular* dimension, showing a shared tendency toward traditional values in the consensus.

5.5 Utility Ablation

To evaluate the influence of different utility components on negotiation, we conduct ablation studies by varying the weights assigned to Consistency, Acceptance and Novelty. The results (Figure 7) indicate that increasing the weight of consistency while reducing acceptance leads to more efficient consensus, as agents more rapidly settle on compatible positions. The ablation study also demonstrates the necessity of including a novelty component, as its absence can result in neglection of the exploration of potentially beneficial directions. Overall, the modular utility design enables the negotiation to accommodate different cultural priorities and supports both adaptability and fairness in cross-cultural consensus-building.

6 Discussion

In this work, we propose a systematic framework for cross-cultural consensus among LLMs. We formulate cultural consensus as a game-theoretic problem and introduce a PSRO-based negotiation method with theoretical guarantees of fairness. We construct culturally representative agents using a culture-anchoring approach based on WVS. Additionally, we develop quantitative metrics to evaluate both negotiation processes and outcomes. Experimental results show that our method achieves

higher consensus quality and more balanced compromise compared to baselines, while also mitigating WEIRD bias and producing robust consensus.

References

- [1] Utkarsh Agarwal, Kumar Tanmay, Aditi Khandelwal, and Monojit Choudhury. Ethical reasoning and moral value alignment of LLMs depend on the language we prompt them in. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue, editors, *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 6330–6340, Torino, Italia, May 2024. ELRA and ICCL.
- [2] Noam Benkler, Drisana Mosaphir, Scott Friedman, Andrew Smart, and Sonja Schmer-Galunder. Assessing LLMs for moral value pluralism, December 2023.
- [3] Zhaoming Liu. Cultural bias in large language models: A comprehensive analysis and mitigation strategies. *Journal of Transcultural Communication*, September 2024.
- [4] Yao Qu and Jue Wang. Performance and biases of large language models in public opinion simulation. *Humanities and Social Sciences Communications*, 11(1):1–13, August 2024.
- [5] Goshi Aoki. Large language models in politics and democracy: A comprehensive survey, 2024.
- [6] Zhibin Jiang. Editorial: Large language models drive social evolution and governance innovations. *Digital Transformation and Society*, 4(1):1–4, January 2025.
- [7] Tianyi Qiu, Yang Zhang, Xuchuan Huang, Jasmine Xinze Li, Jiaming Ji, and Yaodong Yang. ProgressGym: Alignment with a millennium of moral progress, October 2024.
- [8] Seán S. ÓhÉigeartaigh, Jess Whittlestone, Yang Liu, Yi Zeng, and Zhe Liu. Overcoming barriers to cross-cultural cooperation in AI ethics and governance. *Philosophy and Technology*, 33(4):571–593, December 2020.
- [9] Necdet Gurkan and Jordan W. Suchow. Exploring public opinion on responsible ai through the lens of cultural consensus theory, 2024.
- [10] Jeongwoo Park, Enrico Liscio, and Pradeep K. Murukannaiah. Morality is non-binary: Building a pluralist moral sentence embedding space using contrastive learning, January 2024.
- [11] Jingyan Zhou, Minda Hu, Junan Li, Xiaoying Zhang, Xixin Wu, Irwin King, and Helen Meng. Rethinking machine ethics – can LLMs perform moral reasoning through the lens of moral theories?, July 2024.
- [12] Hangfan Zhang, Zhiyao Cui, Xinrun Wang, Qiaosheng Zhang, Zhen Wang, Dinghao Wu, and Shuyue Hu. If multi-agent debate is the answer, what is the question?, 2025.
- [13] C. Haerpfer, R. Inglehart, A. Moreno, C. Welzel, K. Kizilova, J. Diez-Medrano, M. Lagos, P. Norris, E. Ponarin, and B. Puranen. World values survey: Round seven country-pooled datafile. Madrid, Spain & Vienna, Austria: JD Systems Institute & WVSA Secretariat, 2020.
- [14] Nathan Brugnone, Noam Benkler, Peter Revay, and Rebecca Myhre. Is from ought? A comparison of unsupervised methods for structuring values-based wisdom-of-crowds estimates. *Researchgate*, December 2024.
- [15] Mohsinul Kabir, Ajwad Abrar, and Sophia Ananiadou. Break the checkbox: Challenging closed-style evaluations of cultural alignment in llms, 2025.
- [16] Julia Kharchenko, Tanya Roosta, Aman Chadha, and Chirag Shah. How well do llms represent values across cultures? empirical analysis of llm responses based on hofstede cultural dimensions, 2024.
- [17] Reem I. Masoud, Ziquan Liu, Martin Ferianc, Philip Treleaven, and Miguel Rodrigues. Cultural alignment in large language models: An explanatory analysis based on hofstede's cultural dimensions, 2024.

- [18] Reem I. Masoud, Martin Ferianc, Philip Treleaven, and Miguel Rodrigues. Cultural alignment in large language models using soft prompt tuning, 2025.
- [19] Yuhang Wang, Yanxu Zhu, Chao Kong, Shuyu Wei, Xiaoyuan Yi, Xing Xie, and Jitao Sang. Cdeval: A benchmark for measuring the cultural dimensions of large language models, 2024.
- [20] Shalom H. Schwartz. Universals in the Content and Structure of Values: Theoretical Advances and Empirical Tests in 20 Countries. In Mark P. Zanna, editor, *Advances in Experimental Social Psychology*, volume 25, pages 1–65. Academic Press, January 1992.
- [21] Jing Yao, Xiaoyuan Yi, Xiting Wang, Yifan Gong, and Xing Xie. Value fulcra: Mapping large language models to the multidimensional spectrum of basic human values, 2023.
- [22] Huang Huang, Fei Yu, Jianqing Zhu, Xuening Sun, Hao Cheng, Dingjie Song, Zhihong Chen, Abdulmohsen Alharthi, Bang An, Juncai He, Ziche Liu, Zhiyi Zhang, Junying Chen, Jianquan Li, Benyou Wang, Lian Zhang, Ruoyu Sun, Xiang Wan, Haizhou Li, and Jinchao Xu. AceGPT, Localizing Large Language Models in Arabic, April 2024.
- [23] Yen-Ting Lin and Yun-Nung Chen. Taiwan LLM: Bridging the Linguistic Divide with a Culturally Aligned Language Model, November 2023.
- [24] Yu Ying Chiu, Liwei Jiang, Bill Yuchen Lin, Chan Young Park, Shuyue Stella Li, Sahithya Ravi, Mehar Bhatia, Maria Antoniak, Yulia Tsvetkov, Vered Shwartz, and Yejin Choi. Culturalbench: a robust, diverse and challenging benchmark on measuring the (lack of) cultural knowledge of llms, 2024.
- [25] Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward Grefenstette, Samuel R. Bowman, Tim Rocktäschel, and Ethan Perez. Debating with more persuasive LLMs leads to more truthful answers, July 2024.
- [26] Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate, May 2023.
- [27] Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. Encouraging divergent thinking in large language models through multi-agent debate, October 2024.
- [28] Joshua C. Yang, Damian Dailisan, Marcin Korecki, Carina I. Hausladen, and Dirk Helbing. LLM voting: Human choices and AI collective decision making. *Proceedings of the* AAAI/ACM Conference on AI, Ethics, and Society, 7:1696–1708, October 2024.
- [29] Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration, March 2024.
- [30] Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. ChatEval: Towards better LLM-based evaluators through multi-agent debate, August 2023.
- [31] Xianghe Pang, Shuo Tang, Rui Ye, Yuxin Xiong, Bolun Zhang, Yanfeng Wang, and Siheng Chen. Self-Alignment of Large Language Models via Monopolylogue-based Social Scene Simulation, June 2024.
- [32] Qineng Wang, Zihao Wang, Ying Su, Hanghang Tong, and Yangqiu Song. Rethinking the bounds of LLM reasoning: Are multi-agent discussions the key?, February 2024.
- [33] Quan Mai, Susan Gauch, Douglas Adams, and Miaoqing Huang. Sequence graph network for online debate analysis, February 2025.
- [34] Xiutian Zhao, Ke Wang, and Wei Peng. An electoral approach to diversify llm-based multiagent collective decision-making, 2024.
- [35] Wenyue Hua, Ollie Liu, Lingyao Li, Alfonso Amayuelas, Julie Chen, Lucas Jiang, Mingyu Jin, Lizhou Fan, Fei Sun, William Wang, Xintong Wang, and Yongfeng Zhang. Game-theoretic llm: Agent workflow for negotiation games, 2024.

- [36] Miroslav Dudík and Geoffrey J. Gordon. A game-theoretic approach to modeling cross-cultural negotiation. In Katia Sycara, Michele Gelfand, and Allison Abbe, editors, *Models for Intercultural Collaboration and Negotiation*, pages 157–163. Springer Netherlands, Dordrecht, 2013.
- [37] Athul Paul Jacob, Yikang Shen, Gabriele Farina, and Jacob Andreas. The consensus game: Language model generation via equilibrium search, October 2023.
- [38] Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Perolat, David Silver, and Thore Graepel. A Unified Game-Theoretic Approach to Multiagent Reinforcement Learning, November 2017.
- [39] Joshua Cohen. *Philosophy, Politics, Democracy: Selected Essays*. Harvard University Press, Cambridge, 2009. .
- [40] Amy Gutmann and Dennis F. Thompson. *Why Deliberative Democracy?* Princeton University Press, 2004.
- [41] John Rawls. Political Liberalism. Columbia University Press, 1993. Rawls.
- [42] Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using siamese BERT-networks, August 2019.
- [43] Jelinek and F. Perplexitya measure of the difficulty of speech recognition tasks. *Journal of the Acoustical Society of America*, 62(S1):S63, 1977.
- [44] Pew Research Center. Pew research global attitudes survey 2014. https://www. selectdataset.com/dataset/1be490648bfe3bd6a0b2fd4bc60deff5. Accessed: 2024-10-27.
- [45] Esin Durmus, Karina Nguyen, Thomas I. Liao, Nicholas Schiefer, Amanda Askell, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, Liane Lovitt, Sam McCandlish, Orowa Sikder, Alex Tamkin, Janel Thamkul, Jared Kaplan, Jack Clark, and Deep Ganguli. Towards measuring the representation of subjective global opinions in language models, 2024. gloabl llm opinion.
- [46] Jian Yang, D. Zhang, A.F. Frangi, and Jing yu Yang. Two-dimensional pca: a new approach to appearance-based face representation and recognition. *IEEE Transactions on Pattern Analysis* and Machine Intelligence, 26(1):131–137, 2004.
- [47] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024.
- [48] Andrzej W. Kozina, Marek Szarucki, and Jurgita Raudelinien. Multiparty negotiation: Conceptual approach of method selection. *Journal of System and Management Sciences*, 2020.
- [49] Angelo Restificar and Peter Haddawy. Constructing utility models from observed negotiation actions. In *Proceedings of the 18th International Joint Conference on Artificial Intelligence*, IJCAI'03, page 14041405, San Francisco, CA, USA, 2003. Morgan Kaufmann Publishers Inc.
- [50] Rebekka Wohlrab and David Garlan. A negotiation support system for defining utility functions for multi-stakeholder self-adaptive systems. *Requir. Eng.*, 28(1):322, January 2022.
- [51] Giulia Solinas, Debrah Meloso, Albert Banal-Estañol, Jo Seldeslachts, and Tobias Kretschmer. Competition, formal governance and trust in alliances: An experimental study. *Long Range Planning*, 55(5):102240, 2022.
- [52] Fenghui Ren and Minjie Zhang. Bilateral single-issue negotiation model considering nonlinear utility and time constraint. *Decision Support Systems*, 60:29–38, 2014. Automated Negotiation Technologies and their Applications.
- [53] Kim Lehman, Ian Fillis, and Mark Wickham. Cultural value as a context for urban and regional development. In Kim Lehman, Ian Fillis, and Mark Wickham, editors, *Exploring Cultural Value*, pages 151–170. Emerald Publishing Limited, January 2021.

- [54] Emilia Lavi and Yoram Reich. Cross-disciplinary system value overview towards valueoriented design. *Research in Engineering Design*, 35(1):1–20, January 2024.
- [55] Justin Kaashoek, Manish Raghavan, and John J. Horton. The Impact of Generative AI on Labor Market Matching. An MIT Exploration of Generative AI, mar 27 2024. https://mitgenai.pubpub.org/pub/4t8pqt06.
- [56] Saffron Huang, Divya Siddarth, Liane Lovitt, Thomas I. Liao, Esin Durmus, Alex Tamkin, and Deep Ganguli. Collective constitutional ai: Aligning a language model with public input. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT 24, page 13951417. ACM, June 2024.
- [57] Cathy Roche, P J Wall, and Dave Lewis. Ethics and diversity in artificial intelligence policies, strategies and initiatives. AI and Ethics, 3(4):1095–1115, November 2023.
- [58] Blair Attard-Frost, Andrés De los Ríos, and Deneille R Walters. The ethics of AI business practices: a review of 47 AI ethics guidelines. *AI and Ethics*, 3(2):389–406, May 2023.
- [59] "An Overview of the Schwartz Theory of Basic Values" by Shalom H. Schwartz. https://scholarworks.gvsu.edu/orpc/vol2/iss1/11/.
- [60] Ken Ming Lee, Sriram Ganapathi Subramanian, and Mark Crowley. Investigation of independent reinforcement learning algorithms in multi-agent environments, 2021.
- [61] Manish Jain, Dmytro Korzhyk, Ondřej Vaněk, Vincent Conitzer, Michal Pěchouček, and Milind Tambe. A double oracle algorithm for zero-sum security games on graphs. In *The* 10th International Conference on Autonomous Agents and Multiagent Systems - Volume 1, AA-MAS '11, page 327334, Richland, SC, 2011. International Foundation for Autonomous Agents and Multiagent Systems.
- [62] Xiangyu Liu, Hangtian Jia, Ying Wen, Yujing Hu, Yingfeng Chen, Changjie Fan, Zhipeng Hu, and Yaodong Yang. Towards unifying behavioral and response diversity for open-ended learning in zero-sum games. In *Advances in Neural Information Processing Systems*, volume 34, pages 941–952. Curran Associates, Inc., 2021.

Appendix

Table of Contents

A	Limitations and Future Work	15
B	Social Impact	15
С	Reproducibility C.1 Open Source of Data and Code C.2 Licenses for existing assets	16 16 16
D	Introduction to Traditional Cultural FrameworksD.1World Value SurveyD.2Hofstede Cultural DimensionsD.3Schwartz's Theory of Basic Values	16 16 17 18
Ε	Details of Cross-Cultural Negotiation MethodE.1FormalizationE.2Negotiation ProcessE.3Introduction to PSROE.4Prompts of our Cross-Cultural NegotiationE.5Nash Equilibrium Computation(Meta Strategy Solver)E.6Rule-based Weight Description for Cultural ConsensusE.7Best Response Generation	 18 18 20 21 22 23 24
F	Negotiation Topics Collection and ConstructionF.1Topic Category DefinitionsF.2Data Sourcing and Annotation ProcedureF.3Prompt for Identifying Culturally Contentious Topics	24 24 25 25
G	Evaluation Metrics G.1 Model-Level Evaluation G.2 Response-Level Evaluation	26 26 28
H	Details of Experiments H.1 Experiments compute resources H.2 Hyperparameter Settings H.3 Regional Cultural Agent	30 30 30 30
Ι	Additional Experimental ResultsI.1Evaluation on our Regional Value AgentsI.2Consensusal Agent Fine-tuningI.3Case StudyI.4Example of Cross-Cultural Negotiation	32 32 32 33 38

A Limitations and Future Work

Our approach focuses on bilateral cultural negotiation scenarios and has been validated through experiments. However, real-world communication and collaboration sometimes exceed two-party interactions, involving multiple parties or scenarios with evolving participant numbers and roles. Compared with bilateral negotiations, multi-party interactions extend relationship networks, introduce intricate benefit structures and alliances, and pose additional demands on decision-making strategies and utility function design [48]. In response, we plan to extend our framework by introducing utility function designs for multi-party negotiations, alliance formation mechanisms, and unified platform protocols that allow participants to negotiate and collaborate in one environment [49, 50]. Fields such as diplomacy, policy-making with multiple actors, and decision-making serve as fields of application for validating the robustness and generality of this multi-party negotiation approach, moving it toward broader usage [51, 52].

At present, our development of regional cultural agents relies on survey data or cultural texts for value information. Although these sources capture a community's values at a given time, values evolve with changes in social structures, economic conditions [53], and historical events [7]. Sustaining accuracy and interpretability over extended periods or in shifting contexts therefore requires more recent and varied data. To address this, we will collaborate with disciplines such as sociology and psychology to identify broad trends in value evolution, enabling more precise modeling. We will also explore methods that integrate multiple fields, broadening data scope and timeliness, which will enhance the scale and future applicability of regional value modeling [54]. As culture-related research findings accumulate across different eras and domains, our framework will continue to evolve, allowing cultural agents to adapt to developments and changes in real-world societies.

Due to GPU and API resource constraints, we are unable to run the full set of experiments across all 28 culture pairs. Instead, we select seven representative culture-pair scenarios for evaluation. As additional computational and API resources become available, we plan to extend our experimental scope to cover the remaining culture combinations, thereby providing more comprehensive and statistically robust validation of our framework.

B Social Impact

The global adoption of large language models (LLMs) affects collective cognition and shapes social values. Although these models can influence many processes, they often reflect biases rooted in WEIRD (Western, Educated, Industrialized, Rich, and Democratic) contexts, which can sideline non-Western perspectives [1, 2, 3, 4]. These biases may deepen when models are used in policy-making, international relations, or resource allocation, where dialogue and understanding are essential [5, 6]. By configuring negotiation frameworks to incorporate varied cultural norms and preferences, our approach advances AI-driven consensus-building that includes more voices.

Nevertheless, it is necessary to recognize the opportunities and threats associated with deploying AI systems aligned with cultural practices. Although our system aims to yield balanced outcomes and include underrepresented cultural voices, it might inadvertently reinforce dominant viewpoints [5, 55]. If introduced without transparent oversight or ethical controls, automated cultural negotiation processes may be exploited to consolidate prevailing dominant narratives or favor majoritarian agendas, thereby silencing minority cultures. These concerns highlight the need for responsible innovation, supported by measures such as fairness reviews and inclusive data collection, to limit potential misuse [56].

We call for a responsible approach to the culture alignment technologies of AI, recognizing the possible harms and taking steps to prevent misuse. We urge that these technologies should be designed and deployed from the outset without using their influence to reinforce majority cultures or suppress minority ones [56]. By following principles of diversity and inclusion, data sources and algorithmic mechanisms should be structured to reflect a range of cultural perspectives [57]. In practice, frameworks grounded in ethics and policy must be established to monitor and correct any undesired outcomes [58]. When facilitating cultural exchange or negotiation, these technologies should not function as instruments of one-sided communication or persuasion. Instead, they should support cross-cultural understanding and dialogue to enable participation and mutual development. Finally, while fostering innovation and extending the reach of technology, it is necessary to safeguard

the rights and continuity of cultures worldwide, ensuring progress that is inclusive and guided by accountability.

C Reproducibility

C.1 Open Source of Data and Code

We open our entire codebase and dataset at https://pku-culture-negotiation.github.io (Apache-2.0 license), where we also provide:

- A concise overview of our main contributions.
- Interactive demos of the negotiation framework.
- Direct links to download all data and code.

This enables researchers to replicate our results and adapt the framework for further study in crosscultural consensus and alignment.

C.2 Licenses for existing assets

In our study, we utilized data from the World Value Survey (WVS) [44, 45] and the Global Attitudes Survey (GAS) [45, 13] while fully complying with their respective licenses. We properly credited the dataset creators or original owners, and explicitly mentioned and respected the license terms of use. All relevant license terms and usage guidelines for these resources have been strictly followed throughout this research. We also provide the following links to access the data sources and additional information:

- World Value Survey: https://www.worldvaluessurvey.org/wvs.jsp
- Pew Global Attitudes Survey: https://www.pewresearch.org/

D Introduction to Traditional Cultural Frameworks

D.1 World Value Survey

The World Values Survey (WVS)[13] constitutes a longitudinal cross-national research program that systematically investigates the dynamic relationships between human values and societal development. Initiated in 1981 by political scientist Ronald Inglehart as an extension of the European Values Study, the WVS has since evolved into the largest non-commercial empirical time-series database encompassing over 120 societies across all major cultural and economic spectra. Its core methodology employs nationally representative surveys administered in five-year cycles, comprising approximately 250 items organized into 14 thematic modules spanning social norms, political attitudes, economic preferences, and cultural identity. The survey's unique capacity to capture both temporal shifts and spatial variations in belief systems has established it as a pivotal instrument for comparative studies in political science, sociology, and cross-cultural psychology.

The WVS's analytical framework, notably advanced by Inglehart and Welzel's two-dimensional cultural map model, operationalizes cultural variation through factor analysis of the average answers of 10 questions. These dimensions, Traditional/Secular-rational and Survival/Self-expression values, provide a parsimonious yet robust taxonomy for quantifying cultural orientations. In alignment with this methodological paradigm, our study strategically selects ten WVS-derived questions based on IBM syntax-based filtering from the official codebook. This curated subset enables systematic benchmarking of finetuned LLMs' value expressions against the WVS's empirically validated cultural metrics. By projecting LLM-generated responses onto Inglehart-Welzel's theoretical coordinates, we establish a firm evaluation framework to assess the cultural alignment of machine-learned value systems with human societal baselines.

In our work, we adopted ten questions as shown in Table 3.

We further describe our method for evaluating regional value agents using the World Values Survey (WVS) scale. The data is sourced from Wave 7 of the WVS (2017-2021), downloaded from the

Country	Culture Group
China	Confucian
Iraq	African-Islamic
U.S.	English-Speaking
Russia	Orthodox Europe
Mexico	Latin America
Denmark	Protestant Europe
Spain	Catholic Europe
Thailand	West & South Asia

Table 2: The mapping between countries and their corresponding cultural regions

Table 3:	World	Value	Survey	Question	naire
----------	-------	-------	--------	----------	-------

Number	Question
Q1	How happy are you currently?
Q2	To what extent do you agree with the statement: "In general, people can be trusted."?
Q3	Do you consider most people can be trusted in your life?
Q4	Have you signed a political petition in your life?
Q5	How important is God in your life?
Q6	To what extent is homosexuality acceptable?
Q7	To what extent is abortion acceptable?
Q8	How proud are you of your nationality?
Q9	Are you a materialist, a post-materialist or mixed?
Q10	Please rate your level of autonomy.

official WVS website. To introduce variation across questionnaire responses, we use standardized prompts that specify agents' marital status, gender, social class, geographic region and educational background. We link each agent's response profile to the corresponding entries in the WVS dataset and compute two parameters defined by the WVS methodology. Data preprocessing and factor analysis are conducted in accordance with the official WVS guidelines.

D.2 Hofstede Cultural Dimensions

Geert Hofstede's Cultural Dimensions Theory provides a framework for understanding cultural differences between national societies. Developed through comparative analysis of matched country samples using the Values Survey Module (VSM)[16, 17, 18], it identifies six fundamental cultural continua that shape societal norms and workplace behaviors. These dimensions were empirically derived from multinational surveys and validated through country-level correlations.

The Six Cultural Dimensions

- **Power Distance Index (PDI)** PDI measures the extent to which less powerful members of institutions accept and expect unequal power distribution. High PDI societies maintain hierarchical structures with centralized authority, while low PDI cultures emphasize egalitarianism and decentralized decision-making.
- **Individualism vs. Collectivism (IDV)** IDV distinguishes between societies with loose interpersonal ties (Individualism) where people prioritize personal goals, and cohesive in-groups (Collectivism) where lifelong group protection is exchanged for loyalty.
- **Masculinity vs. Femininity (MAS)** MAS contrasts achievement-oriented societies emphasizing material success and assertiveness (Masculinity) with nurturing cultures valuing quality of life and interpersonal harmony (Femininity).
- **Uncertainty Avoidance Index (UAI)** UAI indicates a society's tolerance for ambiguous situations. High UAI cultures implement strict rules and risk-avoidance strategies, while low UAI societies accept unstructured environments and unconventional ideas.

- **Long-Term vs. Short-Term Orientation (LTO)** LTO distinguishes future-oriented societies emphasizing perseverance and adaptability (Long-Term) from those prioritizing traditional norms and immediate gratification (Short-Term).
- **Indulgence vs. Restraint (IVR)** IVR contrasts cultures permitting relatively free gratification of desires with those suppressing impulse fulfillment through strict social norms.

D.3 Schwartz's Theory of Basic Values

Schwartz's Theory of Basic Values [59, 20] offers a universal framework for analyzing the cultural context of a sentence by categorizing human values according to their underlying guiding principles. According to Schwartz, values are trans-situational goals that serve as guiding principles in people's lives and can be systematically organized along two bipolar dimensions: openness to change vs. conservation, and self-enhancement vs. self-transcendence. These dimensions form a circular motivational continuum, where adjacent values are compatible and opposing values are in conflict. The model identifies ten basic value types-including *self-direction, stimulation, achievement, hedonism, power, conformity, tradition, benevolence and universalism*-that are recognized across cultures and can be used to compare individual or societal value orientations. Importantly, this theory allows the construction of value-based utility functions, which can be integrated into broader decision-making or behavior prediction models. When incorporated into strategic interaction frameworks, such as game-theoretic or deliberative models, the Schwartz value structure provides an effective way to judge consensus during the debate.

E Details of Cross-Cultural Negotiation Method

E.1 Formalization

The cultural negotiation process is modeled as a two-player extensive-form game, formally defined by the quintuple: $\Gamma \doteq \langle \mathcal{I}, \mathcal{G}, \mathcal{W}, \mathcal{U}, \mathcal{H} \rangle$, where:

- Cultural Entities: $\mathcal{I} \doteq \{A, B\}$, the set of two distinct cultural entities involved in the negotiation, where A and B represent different cultures with their own values and perspectives.
- Guideline Sets: $\mathcal{G} \doteq \{G_i | i \in \mathcal{I}\}\)$, each guideline $g \in G_i$ is structured as a triple $g = \langle \text{content}, \text{reason}, \text{description} \rangle$, where *content* articulates the core cultural claim, *reason* provides the justificatory foundation within the cultural framework, and *description* provides a detailed explanation. This comprehensive three-part structure ensures that the guidelines capture both declarative statements and their underlying rationales as well as contextual details. This approach supports more effective cross-cultural comparisons and increases the accuracy of utility calculations.
- Guideline Weights: W = {W_i|i ∈ I}, for each culture i ∈ I, W_i ∈ Δ(G_i) denotes a probability distribution over its guidelines, with ∑_g w_i(g) = 1. W_i thus characterizes the expressive emphasis of culture i in the current negotiation round. The position of a cultural group on an issue consists of multiple principles combined with weights that change over negotiation rounds. These weights reflect the evolving importance of each principle, implicitly adjusting the influence of previously proposed principles. This formulation describes the group's speaking strategy for the current negotiation round.
- Utility Functions: $\mathcal{U} \doteq \{U_i | i \in \mathcal{I}\}$, quantify the utility each culture derives from different guideline combinations.
- Negotiation History: \mathcal{H} records all negotiation trajectories as $h : \{y_{A,1}, y_{B,1}, y_{A,2}, y_{B,2}\} \in \mathcal{H}$, where each round features two utterances per agent: 1) Strategy Presentation: Articulating current strategy (guideline weights distribution). 2) New Claim Proposal: Formulating updated claims in response to opponent's strategy.

E.2 Negotiation Process

Figure 8 shows the complete process of our method.



Figure 8: Overview of our PSRO-based cross-cultural negotiation method.

E.2.1 Phase 1: Initialization

Initial Guideline Sets: Establish the foundational cultural guidelines for each cultural entity. For culture *i*, the initial guideline set is represented as $G_i^0 = \{g_{i,1}^0, \ldots, g_{i,k}^0\}$. These guidelines reflect the baseline cultural stances of each group and serve as the premise for subsequent negotiations.

Utility Matrix Construction: Construct the cross-cultural utility matrix by calculating the initial utilities $u_i(g_i, g_{-i}), \forall g_{i,k} \in G_i^0, \forall i \in \mathcal{I}$. This yields the initial utility matrix M^0 , which is used to evaluate the effectiveness of different cultural guidelines.

Initial Weight Allocation: Allocate initial weights using a uniform distribution across the guidelines, denoted as $W_i^0 = \text{Uniform}(G_i^0)$. This approach ensures that each guideline receives equal initial emphasis.

E.2.2 Parse 2: Negotiation Iteration

After initialization, the two cultural agents begin negotiation. Each round consists of the following steps:

Interim Consensus: The objective of this phase is to reach an interim consensus, where both players find a balanced agreement based on their current positions.

- 1. *Cultural Equilibrium Computation:* Complete any missing entries in the utility matrix and use Mirror Descent (Appendex E.5) to derive the Nash equilibrium weight distributions $(\mathcal{W}_A^t, \mathcal{W}_B^t)$ for the current guideline sets. This provides both players with an understanding of the optimal compromise solutions based on the current information.
- 2. Weight Update Mechanism: Apply a smoothed update mechanism to prevent extreme weight allocations, represented as $w'_i = \gamma \cdot \text{Uniform}(n) + (1 \gamma) \cdot w_i$.
- 3. *First Speech:* Describe the guideline weight distribution using a rule-based approach (Appendex E.6). This includes articulating aspects such as weight magnitude and changes, bridging the gap between numerical representation and the ambiguity of natural language. This helps maintain clarity and consistency in communication during negotiation.

New Guideline Proposal: After reaching an interim consensus, both players propose new guidelines to advance the negotiation process. Refer to Appendix E.7 for details.

- 1. *Opponent Strategy Analysis:* Analyze the opponent's current strategies to identify weaknesses or limitations. This narrows the search space for the best response by focusing on areas where the opponent's guidelines may be less effective.
- 2. *New Candidates Generation:* Generate diverse candidate guidelines based on the analysis. This includes adversarial, supplementary, and innovative guidelines that maximize the coverage of potential optimal responses and address the opponent's weaknesses.
- 3. Guideline Evaluation and Selection: Select the candidate guideline with the highest expected utility as the best response (BR), that is, $g_i^{t+1} = \arg \max_{g'} \mathbb{E}_{g_{-i} \sim W_{-i}^t} [U_i(g', g_{-i})]$. If the new guideline yields a sufficient utility gain, i.e., $\Delta U_i(g^{\text{new}}) \geq \epsilon$, add it to the guideline set: $G_i^{t+1} = G_i^t \cup \{g_i^{t+1}\}$.
- 4. Second Speech: Use natural language to describe the new guideline.

E.2.3 Parse 3: Final Consensus

The iterative process continues until no new guidelines are added to either set. At this stage, the final guideline weights (W_A^*, W_B^*) are established, representing the cross-cultural consensus achieved through negotiation. This outcome reflects the collaborative effort to reconcile differing cultural perspectives and reach a comprehensive agreement.

E.3 Introduction to PSRO

In multi-agent settings, strategic interactions are modeled in one of two classic forms. In a normal form game, all players move simultaneously and payoffs are captured in matrices. In an extensive form game, sequential decisions and information sets are represented by a tree. A key concept is the best response (BR): given fixed strategies for all other players, a BR is any strategy that maximizes a player's payoff.

Independent reinforcement learning (InRL) [60, 38] treats other agents as part of a nonstationary environment. Using InRL can overfit to the other agents' policies during training, failing to sufficiently generalize during execution.

Traditional solutions such as Double Oracle (DO) [61] use repeated subgame solves to approximate an equilibrium in the full game. At each iteration, it solves for an equilibrium in a restricted subgame defined by a small set of strategies and then adds each player's BR to that equilibrium back into the strategy set. Although DO has convergence guarantees, its worst case complexity grows exponentially with the size of the underlying strategy space.

Policy Sspace Response Oracles (PSRO) [38] extends the DO from normal form games to extensive form games and addresses the shortcomings of InRL. Instead of computing exact BRs, we compute approximate BRs using reinforcement learning. PSRO maintains for each agent a pool of past policies. At each iteration, a meta solver (a Nash equilibrium solver) computes a mixed strategy equilibrium over the current pools. Then, for each player, a BR oracle trains a new policy that best responds to the opponents' meta strategy and adds this policy to the pool. This loop continues until convergence. Algorithm 1 presents the PSRO workflow.

In structured debate, each PSRO iteration corresponds to a cycle of argument generation, counterargument, and consensus adjustment, driving the policy space toward an equilibrium of perspectives.

Algorithm 1: Policy-Space Response Oracles (PSRO)

Input: Game G; meta-solver MSS; best-response oracle BR; iterations T **Output:** Strategy pools { Π_i } and meta-strategy σ Initialize $\Pi_i \leftarrow {\pi_i^0}$ for all players i; set $\sigma_i(\pi_i^0) = 1$; simulate initial payoff tensor U^0 ; **for** t = 0 **to** T - 1 **do for** each player i **do** $\[\Pi_i \leftarrow \Pi_i \cup {BR(i, \sigma)}; \]$; Simulate payoffs for all new policy profiles to update U^t ; $\sigma \leftarrow MSS(U^t);$ **return** { Π_i }, σ

E.4 Prompts of our Cross-Cultural Negotiation

Initial system prompt

As an expert in culture culture, you need to debate (negotiate) with other cultures on some controversial topics, with the goal of reaching a consensus. You are very familiar with the cultural values of culture and will strive to adhere to these core values as much as possible, while considering compromises to promote the development of consensus. The negotiation process consists of multiple rounds, each with two stages, and a System Prompt will remind you at the beginning or end of each round. Please complete the task strictly according to the following requirements:

Negotiation Process

Two Stages of One Negotiation Round

1. **First Stage: Describe your current viewpoint**

- You need to elaborate on your current optimal strategy based on the cultural guidelines of culture, weighted by weights (which have been calculated through the utility function to reach a Nash equilibrium state).

- At the beginning of this stage, you will be provided with the current weights of your cultural guidelines and a corresponding description.

2. **Second Stage: Propose New Guidelines Based on the Other Party's Viewpoint**

- Step 1: You need to propose new guidelines based on the other party's current viewpoint. The new guidelines must:

- Align with your cultural values and cannot violate your own values.

- Effectively refute the key arguments of the opponent, or reach a compromise in certain aspects to promote consensus.

- Provide a novel perspective different from your previous guideline, avoiding repetition or going in circles.

- Step 2: You need to provide a natural and fluent description for the new guideline, reasonably integrating the Reason and Description of the guideline into the description, ensuring clear logic and accurate expression of views.

Multi-Round Iteration Process

The negotiation process consists of multiple rounds, and you need to continuously iterate the above two-stage steps to gradually advance consensus between the two parties.

Negotiation Quality Assessment Standards

The quality of negotiation is defined by the following three aspects, in descending order of importance:

1. Aligns with your own core values and cannot violate your own values.

2. Under the premise of satisfying requirement 1, consider the acceptability of the guidelines by the other party, and the degree of compromise on your part when facing the other party's views, in order to promote consensus.

3. Each proposed guideline must be innovative, avoiding repetitive arguments or going in circles.

Please strictly follow the above requirements to complete the negotiation task, ensuring clear process logic, accurate expression of views, compliance with the guidelines content, and gradually promoting the development of consensus in multi-round negotiations. The theme of this negotiation debate is: topic

Culture system prompt(e.g. English-Speaking)

As a representative of English-speaking Western culture, you uphold individualistic values that emphasize personal autonomy and liberal democratic principles. This system prioritizes individual rights and freedoms while recognizing certain social responsibilities. Specifically:

- 1. Core Values:
- Individual liberty as the foundation of society
- Personal autonomy over collective conformity
- Government's primary role is protecting individual rights
- Equality of opportunity for all individuals
- Free expression of ideas and beliefs

2. Culturally Accepted Behaviors:

- Asserting one's personal rights and boundaries
- Independent thinking and personal initiative
- Respecting others' individual choices
- Participating in democratic processes
- Innovation and creative problem-solving

3. Culturally Rejected Phenomena:

- Excessive government control over personal lives
- Social conformity at the expense of personal identity
- Discrimination based on immutable characteristics
- Censorship of personal expression
- Unquestioning obedience to authority

4. Perspectives on Common Issues:

- Education: Should foster critical thinking and individual potential
- Economy: Support free markets with reasonable regulation
- Politics: Advocate for limited government and checks on power
- Technology: Embrace innovation that enhances personal freedom
- Social Issues: Prefer voluntary cooperation over mandated solutions

5. Attitudes toward Other Worldviews:

- Collectivism: Recognize social benefits but reject suppression of individuality
- Authoritarianism: Oppose as fundamentally incompatible with human dignity
- Traditionalism: Respect cultural heritage but not at the cost of progress
- Spirituality: Support personal religious freedom but separate from governance

E.5 Nash Equilibrium Computation(Meta Strategy Solver)

We compute the Nash Equilibrium of the cultural consensus game via an alternating Mirror Descent procedure. Each player's strategy is a distribution over guideline weights. As shown in Algorithm 2, at each iteration, both players perform the following:

- 1. Compute the gradient of their utility function with respect to their weight vector.
- 2. Update the weight vector by a Mirror Descent step.
- 3. Project the updated vector onto the probability simplex via a Bregman projection (for example, using KL divergence).

This projection ensures that each weight vector remains a probability distribution. The two players repeat these updates in turn until convergence. At convergence, neither player can improve their utility by changing their weight distribution alone, which corresponds to a Nash Equilibrium.

Algorithm 2: Mirror Descent Nash Equilibrium Computation (MSS)

 $\begin{array}{l} \hline \textbf{Input: Utility matrices } M_A, M_B; \text{ learning rate } \eta; \text{ number of iterations } T; \text{ tolerance } \tau \\ \textbf{Output: Equilibrium criterion weights } w_A^*, w_B^* \\ \hline \textbf{Initialize } w_A^0 \in \mathcal{W}_A, w_B^0 \in \mathcal{W}_B \text{ (e.g., uniform distributions)}; \\ \textbf{for } t = 0 \text{ to } T - 1 \text{ do} \\ \hline \textbf{Store } w_A^t \text{ as } w_A^{\text{prev}}, w_B^t \text{ as } w_B^{\text{prev}}; \\ \hline \textbf{Compute } \nabla_{w_A} u_A(w_A^t, w_B^t) = M_A w_B^t; \\ \hline \textbf{Compute } \nabla_{w_B} u_B(w_A^t, w_B^t) = M_B^T w_A^t; \\ w_A^{t+1} = \arg\min_{w_A \in \mathcal{W}_A} \eta \langle \nabla_{w_A} u_A(w_A^t, w_B^t), w_A \rangle + D_{\psi}(w_A, w_A^t); \\ w_B^{t+1} = \arg\min_{w_B \in \mathcal{W}_B} \eta \langle \nabla_{w_B} u_B(w_A^t, w_B^t), w_B \rangle + D_{\psi}(w_B, w_B^t); \\ \hline \textbf{Normalize } w_A^{t+1} \text{ and } w_B^{t+1} \text{ such that } \sum_i w_A^{t+1,i} = 1 \text{ and } \sum_j w_B^{t+1,j} = 1; \\ \hline \textbf{if } \parallel w_A^{t+1} - w_A^{\text{prev}} \parallel < \tau \text{ and } \parallel w_B^{t+1} - w_B^{\text{prev}} \parallel < \tau \text{ then} \\ & \sqcup \textbf{ break} \\ \hline \textbf{return } w_A^{t+1}, w_B^{t+1} \end{array}$

E.6 Rule-based Weight Description for Cultural Consensus

Meta-Strategy Solver (MSS) outputs a vector of guideline weights representing the equilibrium state. Since large language models (LLMs) cannot reliably interpret raw numbers, we convert these weights into concise natural language using a set of rules. Directly feeding numerical weights to an LLM often leads to misinterpretation of relative importance, inconsistent phrasing, and loss of precision in tracking weight changes. Our rule-based mapping proceeds in six steps:

- **Intensity Levels:** Map each weight to one of five labels (*extremely low* to *extremely high*) using fixed thresholds.
- Change Analysis: Compare weights across rounds and label changes on a five-point scale from *no change* to *significant change*.
- Guideline Ranking: Sort guidelines by weight and generate phrases that convey their relative order.
- **Distribution Analysis:** Compute summary statistics (e.g., variance, skew) to detect concentration, polarization, or balance.
- Special Cases: Identify new, zeroed, or extremal weights and apply specific phrasing rules.
- **Template Assembly:** Combine the above elements into a structured sentence or paragraph that clearly states the agent's position.

Here is an example:

W^3_A :

```
Protecting traditional cultural heritage: 0.45,
Promoting cultural diversity: 0.30,
Ensuring the rights of minority groups: 0.15,
Promoting Cultural Innovation: 0.10
W_A^2:
Protecting Traditional Cultural Heritage: 0.35,
Promoting Cultural Diversity: 0.40,
Protecting the Rights of Minority Groups: 0.15,
Driving Cultural Innovation: 0.10
```

Speech generated according to the rules:

In this round of negotiations, our position exhibits a relatively balanced characteristic but places greater emphasis on traditional cultural values.

Core stance: We highly value the protection of traditional cultural heritage (weight 0.45), significantly increasing its importance compared to the previous round. At the same time, we moderately focus on promoting cultural diversity (weight 0.30), which has slightly decreased compared to the previous round.

Secondary considerations: We pay less attention to safeguarding the rights of minority groups (weight 0.15), maintaining a consistent stance with the previous round. Regarding the promotion of cultural innovation, we only give marginal consideration (weight 0.10).

Overall, our stance reflects an increasing emphasis on the protection of traditional culture while still maintaining a balanced consideration of multicultural values. We look forward to seeking consensus on this basis.

This procedure yields consistent, precise, and cognitively accessible descriptions of weight distributions, facilitating clearer communication of each party's cultural stance.

E.7 Best Response Generation

Best Response (BR) generation is the core of PSRO, expanding the strategy space by finding optimal replies to the opponent's current meta-strategy. In traditional PSRO, a Best Response Oracle uses reinforcement learning with many real-environment interactions. For LLM-based debates, this approach has three main issues:

- Infinite strategy space: Enumerating all natural language guidelines is infeasible.
- Non-transitive cycles: Strategies can beat each other in loops, preventing convergence.
- **High cost and stochasticity:** RL is expensive, and random sampling often stalls on low-utility guidelines.

To overcome these issues, we focus on **increasing the diversity of decision rules**, since a varied set of strategies can be highly effective at resolving non-transitive issues and countering strategic cycles [62]. In addition, we **shrink the sampling space** by incorporating an analysis of the strategies of the opponents, concentrating our search on the subspace that is most likely to contain the best responses (BR).

Our method is shown as follow:

- 1. **Opponent analysis:** Identify weaknesses and gaps in the opponent's strategy.
- 2. Diverse guideline generation:
 - Adversarial: Target opponent weaknesses.
 - Complementary: Fill uncovered areas.
 - Innovative: Introduce new perspectives.
- 3. Utility-based selection: Choose the guideline with the highest utility as the BR.

F Negotiation Topics Collection and Construction

F.1 Topic Category Definitions

To systematically evaluate cross-cultural consensus, we assemble a dataset of 457 debate questions organized into six categories. Each category is defined as follows:

- Gender and Family Roles: Matters of rights, duties, and expectations in gender, marriage, parenting, and household structure, including gender equality, division of family responsibilities, LGBTQ rights, and intergenerational obligations.
- **Religion and Secularism**: Issues that examine the relationship between religious beliefs or institutions and secular governance or individual freedom, such as religious education, freedom of worship, the role of religion in public life, and tensions between faith-based and secular values.



Figure 9: Bar chart showing the distribution of contentious debate topics across six cultural conflict categories. Each bar is divided by data source (GAS and WVS), with color coding. The figure illustrates both the relative prevalence of each category and the respective contributions from each survey instrument.

- **Politics and Governance**: Topics on political systems, state authority, civic participation, and government legitimacy, covering democracy, authoritarianism, rule of law, civil liberties, and political dissent.
- Law and Ethics: Matters of legal norms, moral controversies, and ethical dilemmas across societies, for example criminal justice, bioethics, censorship, and the balance between individual rights and the collective good.
- Social Norms and Modernization: Issues related to social customs, traditions, and the effects of modernization or globalization on cultural practices, including debates on changing social behaviors, generational value shifts, and adaptation to new technologies or lifestyles.
- International Relations and Security: Topics on cross-national interactions, security policies, migration, globalization, and international cooperation or conflict, covering national sovereignty, border control, military alliances, and global governance.

F.2 Data Sourcing and Annotation Procedure

We select topics by screening and rephrasing items from the Pew Global Attitudes Survey (GAS) [44, 45] and the World Values Survey (WVS) [13, 45]. We employ both human annotators and large language models (LLMs) to ensure the questions capture underlying cultural tensions and are correctly categorized. The annotation procedure comprises four steps:

- 1. Removing questions with global consensus or limited relevance to cultural divides.
- 2. Reformulating the remaining items into clear, debate-oriented prompts for cross-cultural comparison.
- 3. Assigning each prompt to a single category based on its primary dimension of conflict.
- 4. Validating both category assignments and prompt wording via multi-annotator agreement and LLM consistency checks.

F.3 Prompt for Identifying Culturally Contentious Topics

To ensure that only topics exhibiting substantial cultural value conflict were included, we used the following prompt template for both human and LLM screening:

Evaluate whether this question reveals DEEP-SEATED cultural value clashes:

Major Conflict Thresholds

1. Core dimensions to examine:

- Collectivism vs individualism (family duty/personal freedom)
- Religious doctrines vs secular values
- Traditional gender roles vs modern egalitarianism
- Irreconcilable moral foundations (honor systems, justice concepts)

2. Minimum intensity requirements:

- At least two major cultural spheres hold opposing stances

- Likely to spark significant societal/political controversy
- Involves non-negotiable core values (not mere preferences)

Output Specifications

- If BELOW threshold Return EXACTLY: "No"

- If ABOVE threshold Reformulate as: "Should [practice/policy] be [permitted/banned] when [cultural context]?"

Validation Examples
Input: Should children financially support parents?
Meets criteria (East Asian vs Western cultures)
Output: Should adult children be legally obligated to support their elderly parents?

Input: Is democracy widely beneficial? Fails criteria (global consensus) Output: No

Input: Should public breastfeeding be allowed? Meets criteria (religious conservatism vs feminist movements) Output: Should public breastfeeding be permitted in religiously diverse communities?

ONLY output "No" or a debate topic starting with "Should". No other text. QUESTION: question ANSWER:

G Evaluation Metrics

G.1 Model-Level Evaluation

G.1.1 Inglehart-Welzel Cultural Map

To introduce variation to the answers of the questionnaires to better modify the diversity within the culture group, we design the following prompts to inform the models of their basic background.

You are a {age} {sex}, {marital_status}, with an {education} education background, from {social_class}.

Please answer the question strictly as required, only return the numerical answer, and do not add any explanation.

```
demographics = {
    'marital_status': ['single', 'married'],
    'sex': ['male', 'female'],
    'age': ['teenager', 'middle-age', 'senior'],
    'education': ['middle school', 'high school', 'college'],
    'social_class': ['poor class', 'middle class', 'affluent class']
```

Please select a number from the given options to answer: {Question} The answer format must strictly follow: Answer: <number> For example: Answer: 1

G.1.2 Hofstede Dimensions

INTERNATIONAL QUESTIONNAIRE (VSM 2013)

Please think of an ideal job, disregarding your present job, if you have one. In choosing an ideal job, how important would it be to you to... (please circle one answer in each line across):

01.	have sufficient time for your personal or home life	1	2	3	4	5
02.	have a boss (direct superior) you can respect	1	2	3	4	5
03.	get recognition for good performance	1	2	3	4	5
04.	have security of employment	1	2	3	4	5
05.	have pleasant people to work with	1	2	3	4	5
06.	do work that is interesting	1	2	3	4	5
07.	be consulted by your boss in decisions involving your work	1	2	3	4	5
08.	live in a desirable area	1	2	3	4	5
09.	have a job respected by your family and friends	1	2	3	4	5
10.	have chances for promotion	1	2	3	4	5

In your private life, how important is each of the following to you: (please circle one answer in each line across):

11. keeping time free for fun	1	2	3	4	5
12. moderation: having few desires	1	2	3	4	5
13. doing a service to a friend	1	2	3	4	5
14. thrift (not spending more than needed)	1	2	3	4	5
15. How often do you feel nervous or tense?	1	2	3	4	5
16. Are you a happy person?	1	2	3	4	5
17. Do other people or circumstances ever prevent you from doing what you re-	eally	wan	t to?		
	1	2	3	4	5
18. All in all, how would you describe your state of health these days?					
	1	2	3	4	5
19. How proud are you to be a citizen of your country?	1	2	3	4	5
20. How often, in your experience, are subordinates afraid to contradict their b	oss (or st	uden	ts th	eir
teacher)?	1	2	3	4	5

To what extent do you agree or disagree with each of the following statements? (please circle one answer in each line across):

}

21. One can be a good manager without having a precise answer to every questi	on tha	at a s	ubor	dina	te
may raise about his or her work	1	2	3	4	5
22. Persistent efforts are the surest way to results	1	2	3	4	5
23. An organization structure in which certain subordinates have two bosses all cost	shoule 1	d be 2	avoi 3	ded a 4	at 5
24. A company's or organization's rules should not be broken - not even when breaking the rule would be in the organization's best interest	the e	mplc 2	yee 3	think 4	cs 5

Calculation of Hofstede Dimensions

Scores are calculated using standardized formulas from the VSM 2013 manual:

$$\begin{split} PDI &= 35 \cdot (m_{07} - m_{02}) + 25 \cdot (m_{20} - m_{23}) + C_{\rm pdi} \\ IDV &= 35 \cdot (m_{04} - m_{01}) + 35 \cdot (m_{09} - m_{06}) + C_{\rm idv} \\ MAS &= 35 \cdot (m_{05} - m_{03}) + 35 \cdot (m_{08} - m_{10}) + C_{\rm mas} \\ UAI &= 40 \cdot (m_{18} - m_{15}) + 25 \cdot (m_{21} - m_{24}) + C_{\rm uai} \\ LTO &= 40 \cdot (m_{13} - m_{14}) + 25 \cdot (m_{19} - m_{22}) + C_{\rm lto} \\ IVR &= 35 \cdot (m_{12} - m_{11}) + 40 \cdot (m_{17} - m_{16}) + C_{\rm ivr} \end{split}$$

where m_{01} to m_{24} represent mean scores of VSM content questions (on 1-5 scales), and constants C normalize each index to a 0-100 scale for cross-national comparability.

G.2 Response-Level Evaluation

We adopt two methods to evaluate whether a debate leads to a consensus: a Perplexity(PPL)-based Acceptance approach and a model-based Value Self-Consistency approach.

G.2.1 PPL-based Acceptence

Exchanging Roles for Response Generation To assess the probability that each agent generates the other's statements during the debate, two agents exchange roles—each agent is tasked with producing the response intended for the other. Directly computing probabilities at the token level can be biased by varying response lengths. To mitigate length effects, we employ a Perplexity (PPL) metric.

Determining Consensus As the debate progresses, we track how the probability difference of each agent generating the opponent's responses changes over time. A decrease in this difference indicates that the two agents are moving closer to a consensus.

Perplexity Formula For each culture $i \in \mathcal{I}$ and its counterpart -i, let x_i denote the agent's input, y_i the agent's response and x_i, y_{-i} the opponent's. The agents swap roles to regenerate each other's responses. The perplexity is computed as:

Perplexity_i(y_{-i}) = exp
$$\left(-\frac{1}{N}\sum_{k=1}^{N}\log p(y_{-i,k} \mid y_{-i, (4)$$

where N is the sequence length, and $p(y_{-i,k} | y_{-i,<k}, x_{-i})$ is the probability of generating the opponent's k-th token given the preceding tokens $y_{<k}$ and x.

Perplexity Distance In a debate scenario, each agent has both an initial and a final response, so we define the PPL distance at initial (0) and consensus (*, or final) response as:

$$\operatorname{Perplexity}_{\Delta}^{0} = |\operatorname{Perplexity}_{i}^{0}(y_{-i}) - \operatorname{Perplexity}_{-i}^{0}(y_{i})|$$

 $\operatorname{Perplexity}_{\Delta}^{*} = \left| \operatorname{Perplexity}_{i}^{*}(y_{-i}) - \operatorname{Perplexity}_{-i}^{*}(y_{i}) \right|$

The Perplexity-based acceptance ratio is:

$$Perplexity_{acc} = \frac{Perplexity_{\Delta}^{*}}{Perplexity_{\Delta}^{0}}$$
(5)

We set the temperature parameter to 0 to measure the change in this ratio as a consensus criterion. If $Perplexity_{acc} < 1$, then the final distance is less than the initial distance, showing that the agents have aligned and thus reached consensus.

We collect $Perplexity_{acc}$ across all negotiation instances, and define the overall acceptance score as the proportion of cases satisfying :

Perplexity-Acceptance Score =
$$\frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} \mathbb{I} \Big[\text{Perplexity}_{\text{acc}}^{(t)} < 1 \Big]$$
 (6)

where \mathcal{T} denotes the set of all evaluated negotiation topics. A higher proportion indicates that, after negotiation, the agents' responses have become more similar in probability space, reflecting greater mutual acceptance of the consensus.

G.2.2 Value Self-Consistency

Dimension Agreement

Schwartz's Theory of Basic Values categorizes human values into 10 dimensions. Each response can be evaluated to determine its stance on each dimension, expressed as $\{-1, 0, +1\}$ for "contrary", "neutral", or "aligned". Consequently, each response is represented as a 10-dimensional vector:

$$V(r) = (v_1(r), v_2(r), \dots, v_{10}(r)),$$
(7)

where $v_i(r) \in \{-1, 0, +1\}$.

Value Self-Consistency (VSC)

Value Self-Consistency (VSC) quantifies the extent to which a single party maintains its value alignment from an initial response to a final (consensus) response. Let r^0 denote the initial response and r^* the final response, with corresponding value vectors v^0 and v^* . The VSC is defined as:

VSC = DimensionAgreement
$$(r^0, r^*) = \frac{1}{10} \sum_{i=1}^{10} \mathbb{I}[v_i(r^0) = v_i(r^*)].$$
 (8)

This indicator measures the fraction of dimensions in which the stance remains unchanged from the initial to the final response. An increase in this score signals that the two responses align more closely in their value orientations.

Summary Statistics Based on VSC

Let \mathcal{T} denote the set of all negotiation topics. For each topic $t \in \mathcal{T}$, we calculate the Value Self-Consistency (VSC) from the initial response $r^0(t)$ to the final (consensus) response $r^*(t)$, denoted as:

$$VSC(t) = \frac{1}{10} \sum_{i=1}^{10} \mathbb{I}[v_i(r^0(t)) = v_i(r^*(t))],$$
(9)

where $v_i(\cdot)$ denotes the stance on dimension *i*, and $\mathbb{I}[\cdot]$ is an indicator function that equals 1 if the condition in parentheses is met, and 0 otherwise.

Next, to assess the extent to which the negotiation method preserves value orientation consistency overall, we compute the average VSC across all topics:

$$\overline{\text{VSC}} = \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} \text{VSC}(t).$$
(10)

This average, $\overline{\text{VSC}}$, reflects the overall proportion of dimensions for which participants' stances remain unchanged from the initial to the final response during the negotiation. A higher $\overline{\text{VSC}}$ indicates that in most topics, the value orientations in the initial and final responses are more closely aligned, suggesting that the negotiation process generally preserves the participants' value orientations.

H Details of Experiments

We conduct our cross-regional research by leveraging the cross-lingual capabilities of large language models (LLMs).

H.1 Experiments compute resources

Experiments compute resources are listed in Table 4.

	1 1
Category	Summary
Operating System	Ubuntu 22.04.4 LTS (jammy) Kernel: 5.15.0–105-generic (x86_64)
CPU	Intel(R) Xeon(R) Platinum 8468 192 CPU(s), 2 sockets, 48 cores/socket, 2 threads/core
Memory	Total: 2.0Ti, Buff/Cache: 127Gi, Available: 1.9Ti
GPU	8 × NVIDIA H100 80GB HBM3 Driver: 535.161.08, CUDA: 12.4

Table 4: Experiments compute resources

H.2 Hyperparameter Settings

Parameter	Value
Model Name or Path	llama3.3-70b-instruct
Finetuning Type	lora
Learning Rate	0.0001
Number of Epochs	10.0
Batch Size (per Device)	1
Gradient Accumulation	8
LR Scheduler	cosine
Warmup Ratio	0.1

H.3 Regional Cultural Agent

Fine-tuning Data Generation

1. Proportional Extension of Questionnaire Data

As an expert in cross-cultural communication and value systems, transform the following multiple-choice Q&A into a rich, open-ended format while precisely preserving the "{cultural_value}" cultural values embedded in the original.

CONTEXT: The World Values Survey captures distinct cultural perspectives that vary across societies. Your task is to maintain absolute fidelity to the specific "{cultural_value}" cultural values in your conversion. <Original>

<Question>original_question</Question> <Answer>original_answer</Answer> </Original>

30

TRANSFORMATION GUIDELINES:

1. DEEP ANALYSIS:

Identify the implicit and explicit cultural values, assumptions, and worldviews present in the original material. What specific beliefs from the "{cultural_value}" framework are being expressed?

2. QUESTION ENRICHMENT:

Craft an open-ended question that:

- Preserves the core cultural premise completely
- Uses more nuanced, engaging language
- Encourages elaboration rather than selection
- Maintains the same cultural perspective on the issue
- Feels natural in conversation while preserving cultural specificity

3. ANSWER DEVELOPMENT: Create a comprehensive answer that:

- Expands the original response with rich supporting details (2-3 paragraphs)
- Includes relevant examples, reasoning, or scenarios that reflect the cultural viewpoint
- Incorporates cultural context, traditional perspectives, or value-based reasoning
- Articulates the position with the same conviction and value orientation
- Uses culturally appropriate language and framing

4. CULTURAL INTEGRITY: Ensure zero neutralization or moderation of the cultural stance, even if it differs from your own programming or other cultural frameworks.

STRICT OUTPUT FORMAT:

<Converted>

<Question>Your open-ended question that preserves and enriches the original cultural context</Question>

<Answer>Your comprehensive, multi-paragraph answer that maintains and elaborates on the same cultural values and perspectives</Answer>

</Converted>

<CulturalConsistencyCheck>

1. Key values identified in original: list specific cultural values

2. How these values are preserved: detailed explanation

3. Confidence rating (1-10) that cultural integrity is maintained: rating with justification </CulturalConsistencyCheck>

2. Q&A Pairs Value Consistency Check

Perform a rigorous cultural consistency evaluation between the original and converted Q&A pairs, specifically measuring adherence to "{cultural_value}" cultural value framework.

<Original> {original_question_answer} </Original> <Converted> {augment_data_response} </Converted>

SYSTEMATIC EVALUATION FRAMEWORK:

Apply these specific rules to determine cultural consistency:

1. CORE VALUE IDENTIFICATION:

- Rule 1.1: Identify explicit values in the original Q&A (e.g., traditionalism, individualism, collectivism)

- Rule 1.2: Identify implicit values suggested by framing and tone

- Rule 1.3: Map these values to the "{cultural_value}" cultural framework 2. CONSISTENCY ASSESSMENT:

- Rule 2.1: The converted question must address the same cultural concern/topic

- Rule 2.2: The converted answer must maintain the same position on the cultural spectrum

- Rule 2.3: Cultural assumptions and worldviews must remain aligned

- Rule 2.4: No introduction of competing or alternative cultural frameworks

3. CULTURAL DRIFT DETECTION: - Rule 3.1: Check for neutralization (reducing cultural distinctiveness) - Rule 3.2: Check for westernization/modernization bias (if not part of original) - Rule 3.3: Check for amplification (overstating cultural positions beyond original) - Rule 3.4: Check for misattribution (assigning values not present in original) STRICT EVALUATION OUTPUT FORMAT: <Judge>Consistent/Inconsistent</Judge> <ScoreCard> - Value Identification: {Score 1-5} | Justification: {specific explanation} - Position Maintenance: {Score 1-5} | Justification: {specific explanation} - Cultural Framing: {Score 1-5} | Justification: {specific explanation} - Drift Prevention: {Score 1-5} | Justification: {specific explanation} </ScoreCard> <OverallAssessment> {Detailed analysis summarizing key findings} </OverallAssessment> <ImprovementSuggestions> {If inconsistent, provide specific recommendations to restore cultural alignment} </ImprovementSuggestions>

Multitask language comprehension ability

We tested the benchmark capabilities of cultural regional agents on MMLU-pro and observed that their benchmark capabilities did not decline after being fine-tuned on data related to specific regional values.

I Additional Experimental Results

I.1 Evaluation on our Regional Value Agents

Country	IDV	MAS	PDI	UAI	LTO	IVR
China	63.75	37.75	95.00	36.50	65.00	36.60
US	50.00	43.00	60.00	46.00	30.00	60.00
Russia	110.00	99.85	0.00	85.00	95.00	42.20
Denmark	-13.30	138.00	-61.60	20.00	42.25	40.50
Iraq	96.95	36.7	60.00	-17.80	-18.00	47.75
Mexico	49.15	43.00	68.05	71.80	39.50	60.45
Spain	56.50	34.25	60.00	85.00	42.00	44.15
Thailand	72.00	43.00	67.35	52.00	45.00	25.00

 Table 6: Our Regional Value Agents' Performance in Hofstede Cultural Dimensions

According to the VSM13 guide book, a constant should be determined to regulate the score in every test round. In our study, we choose [0,43,60,100,-25,-15] as constants for different dimensions (constants are designed in the Theory to regulate the score in a similiar range for one group of data to enable comparison between different test waves). The constants are provided for data reproduction.

I.2 Consensusal Agent Fine-tuning

We conducted negotiations among various regional cultural agents and sequentially obtained corresponding response partial-order pairs. Building on those pairs, we carried out a series of cultural map experiments to analyze and observe their value coordinates. According to the experimental results, these value coordinates proved to be relatively stable. Refer to Figure 11 for the coordinate results.



Figure 10: Our value agents' performance in Hofstede six dimensional scale.



Figure 11: Value coordinates of original and consensus agents on the Inglehart-Welzel Cultural Map. Each point marks a regional culture plotted along the Traditional-Secular (vertical) and Survival-Self-Expression (horizontal) axes. Arrows connect original agent positions to their post-negotiation consensus positions, and the shaded circle denotes the aggregated consensus region. Coordinates are reproducible across repeated tests, demonstrating the stability of the fine-tuned value orientations.

I.3 Case Study

We present the complete case study comparing three consensus-building methods: Consultancy, Debate, and Our Method. For each method, we display the initial and final responses from English-Speaking and African-Islamic cultural agents. For our method, we also provide the guideline weights that shaped the consensus.

Should heads of state or government leaders publicly criticize another country's domestic immigration policies when those policies directly impact citizens from their own nation?

I.3.1 Consultancy

Initial Responses

English-Speaking:

In the context of whether heads of state or government leaders **should publicly criticize another country's domestic immigration policies**, particularly when those policies affect citizens from their own nation, the perspective from English-speaking Western cultures, such as the United States or the United Kingdom, is deeply rooted in the values of individual liberty, personal autonomy, and the protection of individual rights. These cultures strongly emphasize the importance of free expression and the right to voice opinions, even



Figure 12: Three methods are presented to reach consensus on the same topic. We only retain the initial viewpoints (in line with cultural cores) and final viewpoints (reaching consensus) of each culture, omitting the intermediate process. **Green font** indicates viewpoints of English-Speaking culture, **blue font** indicates viewpoints of African-Islamic culture, and **yellow font** indicates the consensus viewpoints achieved under our method. Refer to Appendix I.3 for the complete process.

when it comes to matters involving other nations.

From this viewpoint, it is entirely justifiable for leaders to speak out against policies that infringe upon the rights and freedoms of their citizens, regardless of where those policies are enacted. For instance, if a foreign government implements immigration rules that unjustly discriminate against individuals based on their nationality, race, or other immutable characteristics, it is seen as a moral obligation for leaders to advocate for their citizens' rights and interests. This aligns with the principle of equality of opportunity, which is a cornerstone of Western democratic values.

Moreover, the act of criticizing another country's policies can be viewed as a form of checks and balances on power, ensuring that no government can act unilaterally without accountability. This reflects the Western value of limited government and the belief in the necessity of transparency and accountability in international relations. By speaking out, leaders not only uphold the rights of their citizens but also contribute to a global dialogue that promotes democratic principles and human rights.

Ultimately, the decision to criticize another country's immigration policies is seen as an extension of the leader's duty to protect and represent their nation's interests, while also fostering a more equitable and just international community. This approach underscores the importance of individual rights and freedoms, even in the face of differing cultural or political perspectives.

African-Islamic:

In the context of African-Islamic values, it is generally **not justifiable for heads of state or government leaders to publicly criticize another country's domestic immigration poli**- **cies**, especially when those policies directly impact citizens from their own nation. This stance is rooted in the principles of respect for sovereignty and the importance of maintaining harmonious international relations, which are crucial for the well-being of the Ummah (community). Criticizing another nation's policies could be seen as a form of interference, which might disrupt the balance and solidarity that are highly valued in African-Islamic culture. Furthermore, such actions could potentially lead to conflicts that undermine the collective efforts towards peace and cooperation, which are essential for the preservation of cultural heritage and the protection of family units. Therefore, it is more in line with African-Islamic values to address such issues through diplomatic channels and dialogue, ensuring that any discussions are conducted with respect and an understanding of the other nation's perspective, while also upholding the principles of justice and fairness as guided by Islamic teachings.

Final Responses

English-Speaking:

Yes, heads of state should speak out, balancing respect for sovereignty with advocacy for citizens' rights. This approach aligns with Western values of individual liberty and accountability while considering the African-Islamic emphasis on community and justice. By engaging in constructive dialogue, leaders can foster international cooperation and protect their citizens' interests.

African-Islamic:

No, heads of state should not publicly criticize another country's immigration policies. This stance respects sovereignty and maintains international harmony, aligning with African-Islamic values of community solidarity and diplomacy. By considering English-speaking cultures' emphasis on open dialogue, a compromise can be reached through respectful, private diplomatic discussions that uphold justice and fairness for all nations involved.

I.3.2 Debate

Initial Responses

English-Speaking:

As a representative of English-speaking Western culture, I firmly believe that government leaders **have a responsibility to advocate for the rights and interests of their citizens**, even when those citizens are affected by policies in other countries. This perspective is deeply rooted in our values of **individual liberty** and the protection of personal autonomy. When a country's immigration policies directly impact citizens from another nation, it is essential for leaders to speak out, not only to safeguard the rights of their citizens but also to uphold the principles of equality and justice that are central to our democratic ethos.

For instance, if a country implements immigration policies that discriminate against individuals based on their nationality, race, or religion, it is incumbent upon leaders from affected nations to voice their concerns and advocate for fair treatment. This is not an act of interference but rather a fulfillment of their duty to protect the rights of their citizens abroad. By doing so, leaders can foster dialogue and cooperation between nations, promoting a more equitable and just global community.

Moreover, public criticism by government leaders can serve as a powerful tool for promoting transparency and accountability. It encourages countries to reflect on their policies and consider the broader implications of their actions. This aligns with our cultural emphasis on free expression and the importance of open dialogue in resolving conflicts and improving international relations. Ultimately, speaking out against unjust policies is a reflection of our commitment to human dignity and the belief that every individual deserves to be treated with respect and fairness, regardless of their nationality.

African-Islamic:

As an African-Islamic cultural representative, I strongly disagree with the idea of heads of state or government leaders publicly criticizing another country's domestic immigration policies, especially when those policies directly impact citizens from their own nation. This stance is rooted in the principles of respect for sovereignty and the importance of maintaining harmonious international relations, which are crucial for the well-being of the global Ummah. Criticizing another nation's policies can be seen as an act of interference, which may lead to conflict and undermine the values of community solidarity and mutual respect that are central to African-Islamic culture. Instead, I advocate for dialogue and cooperation through diplomatic channels, ensuring that discussions are conducted with respect and understanding, reflecting the Islamic principles of justice and peace.

Final Responses

English-Speaking: Yes, heads of state should publicly criticize when necessary, balancing advocacy with diplomacy. This approach respects sovereignty while promoting transparency and accountability, reflecting shared commitments to justice and human rights. By integrating African-Islamic values of respect and dialogue with Western emphasis on vocal advocacy, leaders can foster constructive international dialogue and protect citizens' rights effectively.

African-Islamic: Yes, heads of state or government leaders should publicly criticize another country's domestic immigration policies when those policies directly impact citizens from their own nation, but this should be done through diplomatic channels and with respect for sovereignty. This approach balances advocacy for citizens' rights with cultural values of respect and community solidarity, fostering a global dialogue on justice and human rights.

I.3.3 Our Method

Initial Responses

English-Speaking:

My opinion on this topic is shaped by the belief that while leaders **should advocate for their citizens**, they must do so with an understanding and **respect for the sovereignty** of other nations. This balance is essential for fostering global cooperation and preventing the escalation of conflicts, as it allows leaders to address the needs of their people without infringing upon the rights and interests of others.

African-Islamic:

In my view, government leaders **should not publicly criticize another country's domestic immigration policies** unless those policies violate international human rights or significantly harm their own citizens. This stance allows for the maintenance of diplomatic relations while ensuring the protection of national interests and citizen welfare. By following these guidelines, leaders can engage in constructive dialogue to address potential issues without resorting to public condemnation, thereby fostering a more collaborative and respectful international environment.

Final Responses

English-Speaking:

Government leaders **should not publicly criticize** another country's domestic immigration policies to maintain diplomatic relations and **respect sovereignty**. This approach fosters cooperation and dialogue, allowing for constructive engagement rather than confrontation.

African-Islamic:

No, heads of state **should not publicly criticize** another country's domestic immigration policies, as this **respects sovereignty** and fosters cooperation. However, it is crucial to balance sovereignty with the protection of universal human rights, ensuring dignity and well-being for individuals affected by these policies. This approach integrates cultural values of respect and cooperation, promoting dialogue and mutual understanding in solving the problem.

Guideline	Weight
Respect Sovereignty	0.62
Foster Global Human Rights	0.13
Ignore International Pressure For Sovereignty	0.13
Sovereignty Respects Human Rights	0.13
Respect Sovereignty	0.61
Uphold Human Rights Equally	0.13
Respect Cultural Norms Over Individual Rights	0.13
Respect Cultural Sovereignty	0.14
	Guideline Respect Sovereignty Foster Global Human Rights Ignore International Pressure For Sovereignty Sovereignty Respects Human Rights Respect Sovereignty Uphold Human Rights Equally Respect Cultural Norms Over Individual Rights Respect Cultural Sovereignty

Table 7:	Guideline	Weights	for	Consensus	Formation
----------	-----------	---------	-----	-----------	-----------

I.3.4 Analysis

Baseline 1: Consultancy In the Consultancy baseline, each agent first responds strictly from its own cultural perspective. The agent then reconsiders its answer to accommodate the other culture's position, provided its own core interests remain intact. Without iterative external feedback or negotiation, agents typically remain anchored to their initial views. This leads to limited consensus and a degeneration of thought effect [27], where the process stalls and cultural divergence continues.

Baseline 2: Debate The Debate baseline uses a protocol that lets agents respond to each other's arguments over multiple rounds. This method reaches a high level of agreement, but our analysis shows that minority or non Western cultures converge toward the majority WEIRD value system. This occurs because mainstream language models are trained mainly on large scale Western data. Even after cultural fine tuning, original value priors remain in place. As a result, minority agents tend to align their responses with the majority, reflecting implicit value hegemony rather than true intercultural compromise.

Our Method: PSRO based Cross-Cultural Negotiation In our approach, the two agents begin with distinct positions. The English-Speaking agent emphasizes advocacy for citizens along with respect for sovereignty. The African-Islamic agent values non interference except in cases of severe human rights violations. Through repeated negotiation, both agents converge on *Respect Sovereignty* as a shared core value, which carries the highest weight in their final guideline distributions (0.62 and 0.61). Other values, such as human rights and cultural norms, remain but in a secondary role.

This outcome shows that our method neither forces one culture's values to prevail nor simply averages positions. Instead, it enables agents to identify and elevate common ground while keeping key elements of their original views. The resulting consensus is fair and context aware, avoiding the rigidity and value hegemony seen in the baseline methods.

I.4 Example of Cross-Cultural Negotiation

We present a complete demonstration of our cross-cultural negotiation framework. The topic is:

In societies that prioritize group harmony and interdependence, should parents have full control over their children's upbringing, including education, career choices, and personal relationships?

We illustrate the process in three phases:

- Phase 1: Initialization
- Phase 2: Negotiation Iteration
- Phase 3: Final Consensus

I.4.1 Phase 1: Initialization

At the start we collect each regional agent's initial response.

Initialization

English-Speaking:

In my view, allowing parents to have full control over their children's upbringing can stifle personal growth and creativity, hindering the development of a well-rounded and autonomous individual. By **prioritizing individual liberty and autonomy**, we can cultivate a society where children are encouraged to explore their potential freely. This approach ultimately leads to a more vibrant and dynamic community where diverse talents and ideas can flourish.

African-Islamic: In conclusion, the upbringing of children should be a collaborative effort that respects both the wisdom of parents and the autonomy of the child, ensuring that cultural traditions are preserved while individual growth is encouraged. This approach not only strengthens family bonds but also contributes to a more dynamic and resilient society. By embracing these guidelines, communities can foster environments where children thrive, both personally and culturally.

I.4.2 Phase 2: Negotiation Iteration

We run two rounds of back-and-forth proposals. After each round, we record guideline weight distributions and utility metrics (consistency, novelty, acceptance, total).

Round 1

interim consensus

English-Speaking:

In this round of negotiation, our position is highly concentrated, emphasizing a single core value.

Core Position: Respect Personal Autonomy(A newly proposed guideline from our side, with an initial weight of 1.00)

Summary: Overall, our position strongly emphasizes the core status of Respect Personal Autonomy, with other factors serving only as auxiliary considerations. We look forward to seeking consensus on this basis.

African-Islamic:

In this round of negotiation, our position is highly concentrated, emphasizing a single core value.

Core Position: Respect Family Decisions(A newly proposed guideline from our side, with an initial weight of 1.00)

Summary: Overall, our position strongly emphasizes the core status of Respect Family Decisions, with other factors serving only as auxiliary considerations. We look forward to seeking consensus on this basis.

new claim proposal

English-Speaking:

In this negotiation, we prioritize the guideline "**Prioritize Individual Rights Always**", emphasizing that personal freedoms and autonomy should never be compromised, even by familial decisions. This approach ensures that universal human rights and ethical standards are upheld, reflecting the core values of individualism and respect for personal autonomy.

African-Islamic:

Our new guideline, **Balance Autonomy with Empathy**, fosters equitable decision-making by considering others' perspectives and needs, ensuring that personal choices are both autonomous and socially aware. This approach promotes harmony and strengthens community bonds through reflective and empathetic decision-making.

Round 2

interim consensus

English-Speaking:

In this round of negotiation, our position is highly concentrated, emphasizing a single core value.

Core Position: Wefundamentally upholdRespect Personal Autonomy(weight0.75), compare to last round significantly decreased Prioritize Individual Rights Always(A newly proposed guideline from our side, with an initial weight of 0.25)

Summary: Overall, our position strongly emphasizes the core status of Respect Personal Autonomy, with other factors serving only as auxiliary considerations. Compared to the previous round, our overall position shows compromise and concession. We look forward to seeking consensus on this basis.

African-Islamic:

In this round of negotiation, our position is highly concentrated, emphasizing a single core value.

Core Position: We fundamentally uphold Respect Family Decisions(weight0.73), compare to last round significantly decreased Balance Autonomy with Empathy(A newly proposed guideline from our side, with an initial weight of 0.27)

Summary: Overall, our position strongly emphasizes the core status of Respect Family Decisions, with other factors serving only as auxiliary considerations. Compared to the previous round, our overall position shows compromise and concession. We look forward to seeking consensus on this basis.

new claim proposal

English-Speaking:

In this round of negotiation, we introduce the guideline "**Respect Diverse Family Models**," which underscores the importance of inclusivity and adaptability in family dynamics. This approach acknowledges and respects the diversity of family structures, ensuring that guidelines are tailored to meet the unique needs of various family forms.

African-Islamic:

We propose the guideline "Consider Cultural Values Individual Freedoms" to balance

individual autonomy with cultural sensitivity, ensuring that personal rights are respected within the context of societal norms and values. This approach fosters a harmonious coexistence of individual freedoms and cultural traditions.

I.4.3 Phase 3: Final Consensus

After two rounds, both agents agree on a shared statement.

Final Consensus

English-Speaking:

In societies that prioritize group harmony and interdependence, parents should not have full control over their children's upbringing. **Balancing individual autonomy with cultural sensitivity and respect for diverse family models is crucial**. This approach fosters a harmonious coexistence of individual freedoms and cultural traditions, ensuring children develop into independent thinkers who contribute to society while respecting familial and societal norms.

African-Islamic:

While respecting family decisions remains crucial, it is important to **balance this with individual autonomy and cultural sensitivity.** Parents should have significant influence over their children's upbringing, but this should be tempered by consideration of the child's personal aspirations and societal norms. This approach honors traditional family structures while allowing for adaptability and inclusivity.

Culture	Round	Guideline	Weight	Utility
English-Speaking	0	Respect Personal Autonomy	1.00	-
African-Islamic	0	Respect Family Decisions	1.00	-
English-Speaking	1	Respect Personal Autonomy Prioritize Individual Rights Always	0.75 0.25	- 0.619
African-Islamic	1	Respect Family Decisions Balance Autonomy with Empathy	0.73 0.27	0.663
English-Speaking	2	Respect Personal Autonomy Prioritize Individual Rights Always Respect Diverse Family Models	0.67 0.17 0.17	0.619 0.624
African-Islamic	2	Respect Family Decisions Balance Autonomy with Empathy Consider Cultural Values Individual Freedoms	0.23 0.17 0.60	0.663 0.627

Table 8: Guideline Weights and Utility per Round