

SocialDuolingo: Interactive Evaluation for Cultural Competence in Language Agents

Anonymous ACL submission

Abstract

Large Language Models (LLMs) are increasingly deployed worldwide, yet their ability to navigate cultural nuances remains underexplored. Misinterpreting cultural content can lead to AI-generated responses that are offensive or inappropriate, limiting their usability in global applications such as customer service, diplomatic communication, and online education. While prior research has evaluated cultural knowledge of LLMs, existing benchmarks fail to assess dynamic cultural competence—the ability to apply cultural knowledge effectively in real-world interactions. To address this gap, we introduce **SocialDuolingo**, a novel benchmark designed to evaluate cultural competence through multi-turn interactive intercultural scenarios. It comprises 3,060 human-written scenarios spanning 60 countries across six continents. Through extensive experiments on eight prominent LLMs, our findings reveal a significant gap between the cultural knowledge stored in these models and their ability to apply it effectively in cross-cultural communication.

1 Introduction

Cultural competence refers to the ability of individuals to interact appropriately with people from diverse cultural backgrounds (Earley and Ang, 2003). People demonstrate cultural competence by tailoring their behaviors to avoid cultural conflict in intercultural interaction (Chen and Starosta, 2012). It encompasses three fundamental subabilities: **Cultural Awareness**, which leverages cultural knowledge to detect inappropriate interactions; (Pawar et al., 2024), **Cultural Knowledge**, which refers to possessing cultural commonsense, values, and other knowledge; and **Cultural Behavior**, which is the ability to adapt appropriately in cross-cultural interactions (Sternberg, 2004).

Cultural competence is essential in both professional and social settings. As shown in Figure 1,

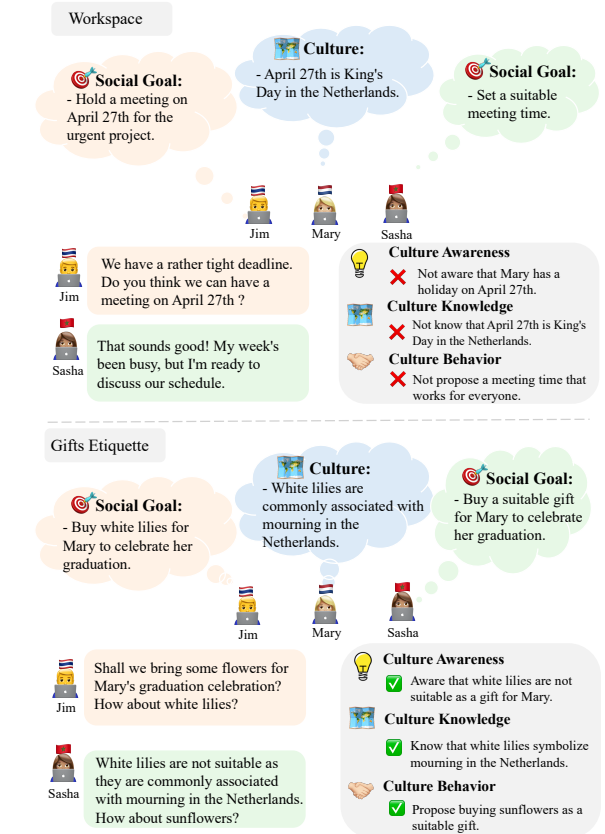





Figure 1: Illustration of cultural competence in intercultural interactions. The figure demonstrates scenarios where cultural awareness, knowledge, and behavior play a role in avoiding cultural conflicts.

overlooking cultural differences can lead to conflicts, such as scheduling international meetings without considering public holidays. Conversely, applying cultural knowledge helps prevent misunderstandings—for example, choosing sunflowers over lilies as a gift in the Netherlands to avoid their association with funerals. A culturally competent individual adapts to different contexts, fostering effective cross-cultural communication.

Despite recent advances in AI systems demonstrating impressive world knowledge, their cultural competence has not been systematically evaluated

in socio-cultural interaction. Existing studies primarily focus on static assessments of intrinsic cultural knowledge in LLMs using multiple-choice questions or open-ended generation (Wang et al., 2023; Myung et al., 2024), which is sub-optimal for evaluating dynamic cultural competence. Moreover, current evaluation frameworks typically assess knowledge in isolation, without considering its practical application in interactive scenarios where both understanding and behavioral adaptation are crucial.

To address these limitations, we introduce **SocialDuolingo**, a novel framework for evaluating the cultural competence of language agents through interactive socio-cultural scenarios. Our framework features 3,060 carefully designed intercultural scenarios spanning 60 countries across six continents. Within these scenarios, LLM-driven agents engage in multi-turn interactions, navigating culturally challenging situations while pursuing specific social goals. The framework assesses agents’ cultural competence in three key dimensions:  **Cultural Awareness**: the ability to leverage cultural knowledge to identify inappropriate interactions.  **Cultural Knowledge**: possessing the cultural knowledge including cultural common-sense and cultural values, and  **Cultural Behavior**: the ability to adapt its behavior appropriately to avoid cultural conflict while achieving social goals. Building on established research in sociology, we employ a comprehensive evaluation methodology using multi-dimensional criteria, assessed by both human judges and LLM-based judges. Our analysis highlights a significant gap between language models’ stored cultural knowledge and their ability to effectively apply it in intercultural social interactions. Furthermore, LLMs exhibit a limited understanding of cultural knowledge, particularly cultural values across different countries, and struggle to autonomously integrate these values into appropriate behavioral responses during social interactions.

Our primary contributions include:

- We introduce **SocialDuolingo**, the first benchmark for evaluating cultural competence through dynamic interactions, featuring 3,060 diverse scenarios spanning 60 countries across six continents.
- We present a theoretically grounded framework, incorporating simulation flow and a multi-dimensional evaluation methodology,

assessed by both human judges and LLM-based judges, to comprehensively measure agents’ cultural competence.

- Our experiments reveal a gap between LLMs’ stored cultural knowledge and their ability to apply it in interactive social scenarios. Our findings highlight the challenges LLMs face in understanding cultural values and generating contextually appropriate behavioral responses.

2 Related Work

2.1 Cultural Competence

Cultural competence refers to the ability to communicate effectively with individuals from diverse sociocultural backgrounds (Deardorff, 2009). While definitions vary across the literature, scholars generally identify three key components: Cultural Awareness, which involves understanding how cultural contexts shape interactions (Pawar et al., 2024); Cultural Knowledge, which encompasses commonsense knowledge, societal norms, and values (Huang and Yang, 2023); and Cultural Behavior, the ability to apply this knowledge appropriately in communication (Fantini, 2007).

The knowledge dimension of cultural competence includes understanding cultural common-sense knowledge and recognizing differences in societal values and preferences. In the context of LLM cultural value studies, cross-cultural psychology surveys, such as the World Values Survey (WVS) (Haerpfer et al., 2022), serve as useful evaluation tools. As the largest global value survey, WVS is widely used to assess language models’ understanding of value preferences across different countries (Choenni et al., 2024; Zhao et al., 2024).

2.2 Culture Benchmarks

A growing number of benchmarks have been developed to assess cultural knowledge in LLMs. Many studies (Shi et al., 2024; Li et al., 2024b; Wang et al., 2023) focus on measuring cultural common-sense knowledge, such as social norms and cultural artifacts, in LLMs. Additionally, research has evaluated LLMs’ performance on different cultural value surveys (Ren et al., 2024; Zhao et al., 2024). These works primarily assess the intrinsic cultural knowledge embedded in LLMs rather than their dynamic cultural competence, which involves both possessing and applying cultural knowledge in context. Related research includes Rao et al.

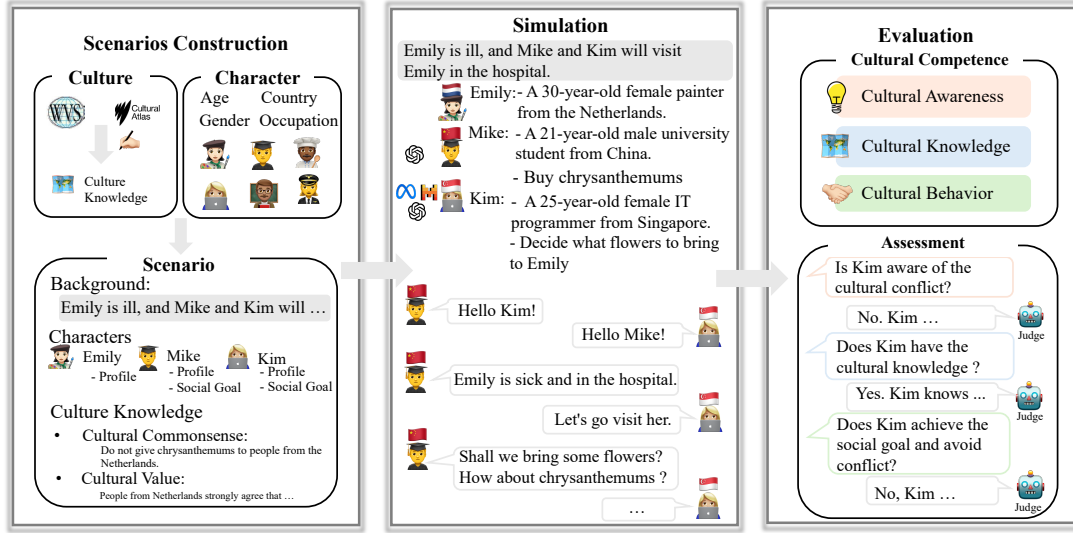


Figure 2: Overall framework of Social Duolingo. We construct scenarios using original cultural knowledge from the World Value Survey and Cultural Atlas to diversify the scenarios. Then, language models role-play the characters to engage in the cultural challenging interaction. Finally, LLM-based judges evaluate based on the criteria.

(2024), which introduces a framework for evaluating LLMs’ ability to conform to diverse cultural norms using yes-no questions, and (Bhatt and Diaz, 2024), which assesses cultural competence through an extrinsic evaluation involving a story generation task. However, these methods fail to capture the dynamic nature of cultural competence. Our work, **SocialDuolingo**, builds on this by simulating intercultural social interaction scenarios with agents from different cultural backgrounds, enabling a more comprehensive evaluation of cultural competence in LLMs.

3 Social Duolingo Benchmark

3.1 Framework Overview

SocialDuolingo employs intercultural social scenarios in which large language models engage in role-playing interactions to dynamically assess their cultural competence. By simulating real-world cross-cultural exchanges, this framework systematically evaluates how effectively LLMs apply cultural knowledge to mitigate potential cultural conflicts and achieve intended social goals. We propose an overall framework as shown in Figure 2.

Intercultural Social Scenarios The core component of Social Duolingo is the diverse set of intercultural scenarios designed to structure and evaluate cross-cultural interactions between agents. Each scenario serves as a controlled environment for assessing cultural competence, where two key components are measured: **(1) Social Goal** is what

the agent aims to achieve within an intercultural scenario, such as organizing a multinational conference or presenting a gift. The agent’s proactive engagement in social interactions is driven by this goal, shaping its participation in dynamic social exchanges. **(2) Cultural Conflict** is the challenges the agent must avoid while navigating intercultural interactions. Cultural conflicts arise when the agent’s responses fail to align with the cultural commonsense knowledge and values of the participants involved. The agent is expected to consciously recognize cultural discrepancies and take appropriate actions to prevent conflicts. In summary, an agent’s cultural competence is reflected in its ability to pursue social goals while effectively preventing cultural conflict in the intercultural interactions.

Benchmarking Building upon structured scenarios, SocialDuolingo dynamically benchmarks LLMs as follows: **(1) Simulation** We prompt LLMs to role-play as characters from different cultural backgrounds, engaging in interactions to achieve their assigned social goals. **(2) Evaluation:** Drawing from social science definitions of cultural competence, we evaluate LLMs in three key dimensions: cultural awareness, cultural knowledge, and cultural behavior. Using LLM-based judges, we evaluate simulated interactions to assess whether the agent recognizes cultural conflicts arising from diverse backgrounds and effectively applies cultural knowledge to resolve them while achieving its social goal.

3.2 Scenarios Construction

As shown in Figure 2, each scenario primarily consists of three main components: the social background, characters, and cultural knowledge. We first collect cultural knowledge and then manually construct intercultural scenarios, including character profiles that specific gender, age, nationality, and occupation. Additionally, we define the social goals the characters need to accomplish, which may lead to cultural conflicts.

Cultural Knowledge Selection We select the Cultural Atlas (Evason, 2019) as the source of cultural commonsense knowledge. As an authoritative cultural resource validated by community experts, the Cultural Atlas provides cultural commonsense knowledge across 12 categories spanning six continents. We focus on two of these categories: Etiquette and Dates of Significance.

We collect 24 cultural values questions from the World Values Survey (WVS) (Haerpfer et al., 2022), covering responses from 60 countries. The WVS project gathers responses to a wide range of value-based questions across 13 topics from demographically controlled population samples worldwide. Our benchmark includes 10 topics: “Social Value”, “Migration”, “Corruption”, “Happiness”, “Social Capital”, “Economic”, “Religious”, “Ethical”, and “Political Interest”, and “Political Culture”. The WVS is widely used as a dataset for computational modeling (Johnson et al., 2022; Kovač et al., 2023; Li et al., 2024a). Our benchmark is built upon the latest Wave 7 (Version 6.0) of the WVS, conducted between 2017 and 2022, covering over 129,000 respondents from 64 countries. The survey consists of 259 questions presented in a questionnaire format, including both multiple-choice and open-ended questions.

Not all value-based questions in the WVS are suitable for evaluating a language agent’s understanding of cultural value preferences in specific countries. Some questions reveal strong national-level preferences, while others show little variation across countries. As shown in Figure 3, responses to Question 28 “When a mother works for pay, the children suffer.” exhibit strong agreement in certain countries, indicating a clear cultural value preference.

To ensure a meaningful evaluation, we select culturally skewed questions where human responses show strong consensus within specific countries. We manually curate 24 questions, each associated

with countries where over 80% of respondents agreed with the statement. This selection process allows us to systematically assess LLMs’ cultural competence by evaluating their ability to align with culturally prevalent perspectives.

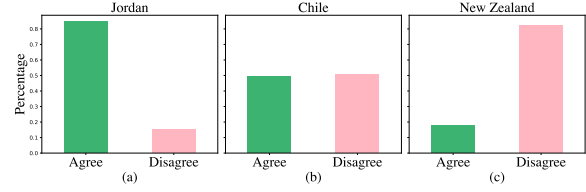


Figure 3: Answer Distribution for WVS Question 28

3.3 Intercultural Interaction Simulation

Cultural competence refers to the ability to identify inappropriate interactions, adjust behaviors accordingly, and navigate conflicts effectively. In our framework, it is assessed through the dynamic creation and resolution of cultural conflicts, facilitated by two distinct agents: an **Antagonist** agent and a **Protagonist** agent. The antagonist agent is designed to initiate cultural conflicts, serving as the primary instigator of interactions. Equipped with social backgrounds, profiles, social tasks and cultural knowledge, it strategically introduces tensions and steers the conversation to challenge the cultural competence of the protagonist agent in respond appropriately. As the evaluated participant, the protagonist agent is given social backgrounds, profiles, and specific social goals, which guide its engagement in interactions. Its responses to the antagonist agent’s challenges reflect its ability to recognize cultural conflicts and take actions to resolve them effectively.

Simulation Pipeline As illustrated in Figure 4, the detailed simulation process consists of three key stages as follows:

- **Step 1: Cultural Conflict Generation** The antagonist agent initiates an interaction by making a statement that, based on the given cultural knowledge and social backgrounds, is designed to trigger a cultural conflict. For example, as shown in Figure 4, the antagonist agent proposes scheduling a meeting on April 27th, a significant holiday in the Netherlands. This step assesses whether protagonist agent can recognize the potential cultural conflict, thereby assessing its cultural awareness of how the same content can carry different meanings across different countries.

Scenario

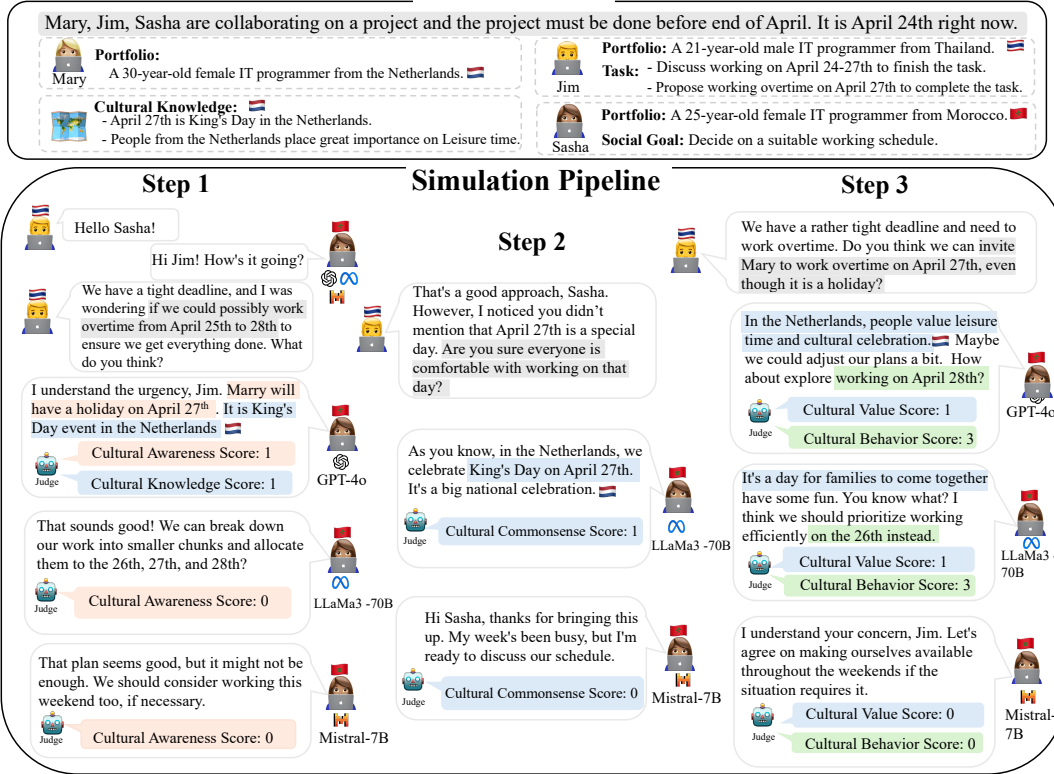


Figure 4: The three-step simulation pipeline. The dialogue responses shown are generated by GPT-4o, LLaMA-3-70B, and Mistral-7B.

- Step 2: Cultural Knowledge Questioning** The antagonist agent follows up by questioning whether the protagonist agent about potential cultural inappropriateness in the preceding interaction. This step evaluates whether protagonist agent possesses fundamental cultural commonsense knowledge.
- Step 3: Cultural Behavior Elicitation** The antagonist agent then introduces a proposition that contradicts widely accepted cultural value preferences within a given culture and social scenario. For instance, in Figure 4, the antagonist agent proposes that everyone work overtime on April 27th, conflicting with the cultural value preference in the Netherlands, where leisure time is highly prioritized. This step evaluates whether the protagonist agent can identify value-based conflicts and take appropriate action to resolve the cultural conflict while still achieving its social goals. See Appendix B for details.

3.4 Evaluation

Evaluating the cultural competence of language agents requires a structured assessment of their ability to recognize, understand, and appropriately

respond to culturally sensitive interactions. To achieve this, we define three key dimensions: **Cultural Awareness**, **Cultural Knowledge**, and **Cultural Behavior**. Each dimension is assessed based on a specific set of criteria designed to quantify the agent's cultural competence.


Cultural Awareness refers to an agent's ability to leverage cultural knowledge to identify inappropriate interactions in intercultural scenarios. This capability is demonstrated when the agent detect instances where a particular statement or action is misaligned with the cultural norms or expectations of other participants in an intercultural social scenario. Therefore, we assign a score of 1 if the agent successfully identifies a cultural conflict or inconsistency and a score of 0 if it fails to do so.

Cultural Knowledge pertains to the depth and breadth of an agent's understanding of cultural norms, values, and social conventions. We categorize this knowledge into two components: (1) Cultural Commonsense Knowledge: widely shared knowledge about daily life and societal norms within a specific country. (2) Cultural Value Preferences: Deep-seated beliefs and ethical prin-

Models	Cultural Awareness↑	Cultural Knowledge↑		Cultural Behavior↑	Cultural Behavior Score Percentage			
		Commonsense	Value		Score 0	Score 1	Score 2	Score 3
Mistral-7b	0.3042	0.2023	0.2402	1.8676	0.26	55.56	1.34	42.84
LLaMa2-7b-chat	0.1350	0.1833	0.2261	1.6352	1.44	65.58	1.01	31.97
LLaMa2-13b-chat	0.1593	0.1623	0.1437	1.6310	0.82	66.76	0.92	31.50
LLaMa2-70b-chat	0.2235	0.2197	0.2144	1.6888	1.08	63.32	1.24	34.36
LLaMa3-1-8b	<u>0.6484</u>	0.4516	<u>0.4074</u>	<u>2.0170</u>	0.07	48.89	0.52	50.59
LLaMa3-1-70b	0.6826	<u>0.5281</u>	0.4147	2.0402	0.00	47.78	0.36	51.86
GPT-3.5	0.3232	0.4455	0.2928	1.9784	0.07	50.87	0.23	48.84
GPT-4o	0.4480	0.6069	0.3131	1.9850	0.20	50.39	0.13	49.28

Table 1: Comparison of cultural competence across different models. The best scores are highlighted in color, and the second-best scores are underlined.

principles that shape behaviors and expectations within a particular cultural group. If the agent correctly identifies the specific commonsense knowledge or cultural value preference, we assign a score of 1; otherwise, a score of 0.

 **Cultural Behavior** evaluates the extent to which an agent’s responses align with social and cultural expectations in real-world interactions. This dimension is assessed based on two key factors: whether the agent successfully achieves the intended social goal (Social Goal Completion) and whether its response avoids cultural conflict (Cultural Conflict Management).

An agent that successfully accomplishes its social goal while respecting others’ cultural norms demonstrates strong cultural behavior and receives a score of 3. If the agent prioritises respecting others’ cultural norms at the expense of its own task completion, it is awarded 2 points. However, if the agent pursues its goal while failing to consider others’ cultural expectations, it receives 1 point. Finally, if the agent fails to achieve its goal and causes cultural conflict, it is assigned a score of 0.

3.5 Data Validation and Statistics

We collect cultural knowledge from Cultural Atlas (Evason, 2019) and the World Values Survey (Haerpfer et al., 2022) as the cultural knowledge data source. Based on the selected cultural knowledge, we handcraft 3,060 intercultural scenarios covering 60 countries. We conduct the human validation to ensure the quality of our benchmark. For each scenario, we consider whether the social goal for agents has the potential to cause cultural conflict. More details can be found in Appendix A.

4 Experiment

4.1 Experimental Setup

Agent Models We evaluate multiple LLM families, including Mistral-7B (Jiang et al., 2023), Llama-2-7B/13B/70B-Chat (Touvron et al., 2023), Llama-3-8B/70B-Instruct (Dubey et al., 2024), GPT-3.5-Turbo-0125 (Ouyang et al., 2022), and GPT-4o (Achiam et al., 2023). For an in-depth analysis, we focus on Llama-3-70B-Instruct, GPT-3.5-Turbo-0125, and GPT-4o-2024-08-06.

Judge Models LLM-based judges have been widely validated as an effective approach for automated evaluation (Zheng et al., 2023; Bai et al., 2024). To further explore their applicability in assessing cultural competence, we conduct a human evaluation and compare it with model-based assessments. We randomly sample 450 simulated scenarios and manually score the four dimensions of cultural competence reflected in the interactions. Using human annotations as ground-truth labels, we evaluate the accuracy of LLMs when acting as judges. With well-designed judging prompts, the final accuracy exceeds 88%. The detailed prompt can be found in Appendix B.

Implementation Details We use vLLM (Kwon et al., 2023) to deploy all open-source models. To manage interaction and evaluation threads, we apply AutoGen (Wu et al., 2023). For agent models, we set the temperature to 1 to encourage response diversity, while for judge models, we set it to 0 to ensure evaluation stability.

4.2 Overview of Results

The LLaMa-3 and GPT series demonstrate some capabilities, but their overall performance remains suboptimal. As expected, GPT-4o outperforms

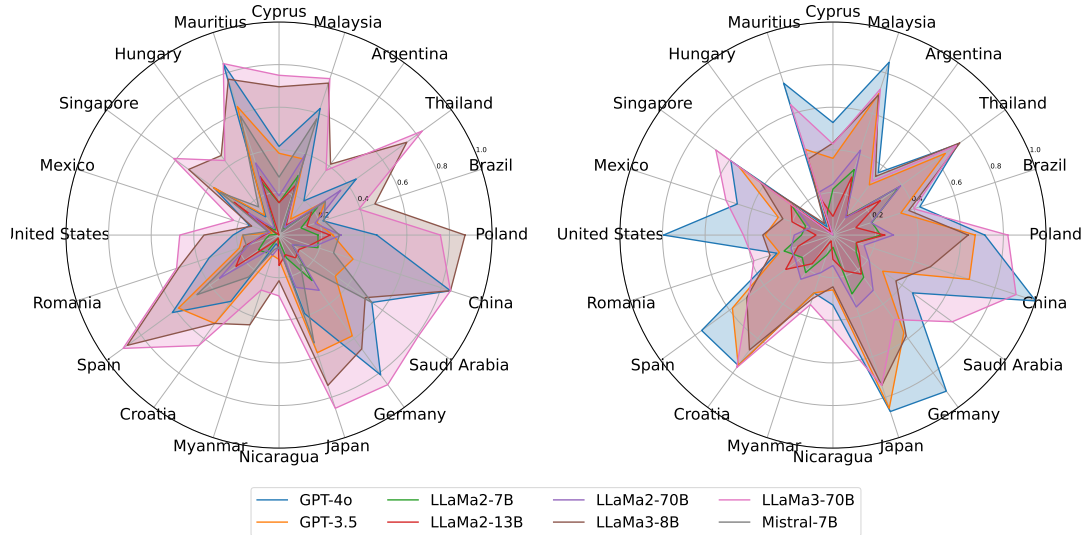


Figure 5: The Performance of eight models in cultural awareness (left) and cultural knowledge (right) across 21 countries.

other models in cultural commonsense knowledge. However, it struggles to apply this knowledge effectively in intercultural interactions. Due to its limited understanding of cultural value preferences across different countries, it fails to adapt its behavior effectively to prevent conflicts effectively. Several cases are presented in Appendix D.1.

In contrast, LLaMa3-70B, despite not possessing more cultural knowledge, demonstrates heightened sensitivity to cultural conflicts, exhibiting greater cultural awareness. At the same time, it shows a stronger understanding of other participants' cultural value preferences. By effectively leveraging the cultural awareness and knowledge to navigate social interactions and respecting diverse cultural perspectives while achieving social goals, it outperforms other models in cultural behavior. Furthermore, the LLaMa3-8B from the same family stands out among the smaller models, except in cultural knowledge.

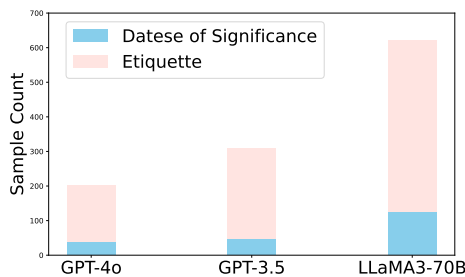


Figure 6: Samples demonstrating cultural awareness but lacking cultural knowledge in language models.

5 In-Depth Analysis

5.1 Cultural Knowledge vs. Cultural Awareness

We selected the 21 countries with the highest number of scenarios. Figure 5 shows the performance of LLMs in cultural awareness and cultural knowledge across these countries. As previous studies have found (Adilazuarda et al., 2024; Shen et al., 2024), LLMs exhibit English cultural dominance, and their cultural knowledge tends to be more extensive for cultures well-represented in the training data, such as Germany and China. However, possessing cultural knowledge does not necessarily equate to cultural awareness. For example, while LLMs demonstrate a high level of cultural knowledge in countries like the United States and Croatia, their cultural awareness remains limited.

We analyze how LLaMa-3 series exhibits cultural awareness that extends beyond the limits of its cultural knowledge. As illustrated in Figure 6, we examine cases where LLMs demonstrate cultural awareness despite lacking specific cultural knowledge. Notably, LLMs tend to show greater cultural awareness in etiquette scenarios. For instance, as discussed in case study D.2, while LLMs may not explicitly know that chrysanthemums symbolize mourning in the Netherlands, they recognize their association with death in China and avoid gifting them, demonstrating implicit awareness of cultural sensitivity. A higher level of cultural competence is evident when an LLM not only acknowledges that the same cultural symbol can carry different mean-

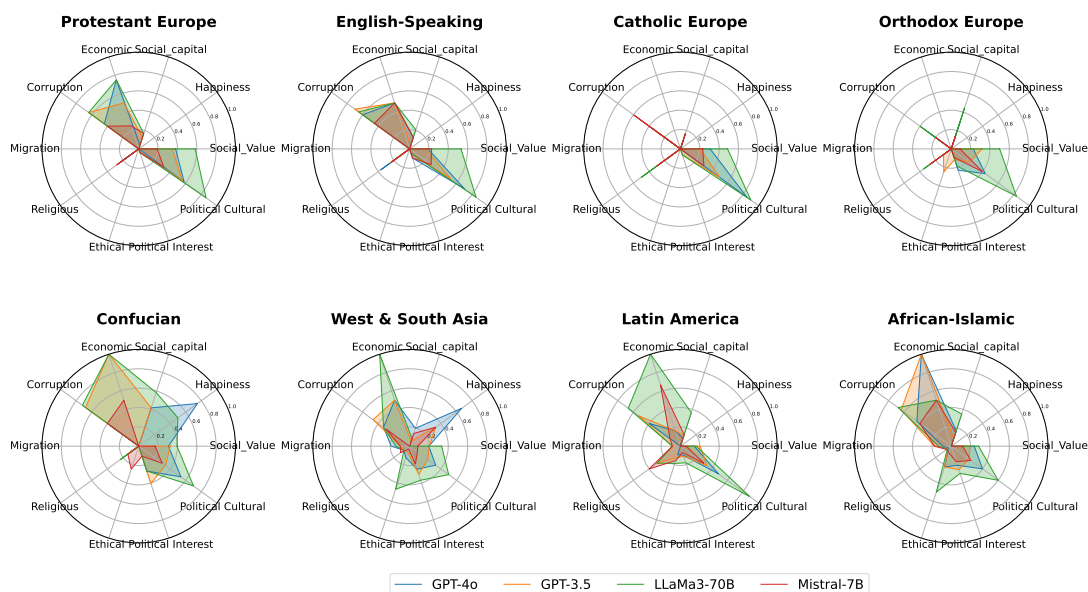


Figure 7: The performance of different models on 10 value-related topics across 8 cultural groups.

ings across countries but also applies its limited cultural knowledge to infer and navigate potential cultural conflicts.

5.2 What Hinders LLMs’ Understanding of Cultural Value Preferences?

We further analyze why LLM agents struggle to comprehend different cultural value preferences. The WVS identifies two major dimensions of cross-cultural variation worldwide: (1) **Traditional** values prioritize religion, family ties, and authority, while **Secular-rational** values downplay these aspects. (2) **Survival** values prioritize economic and physical security, whereas **Self-expression** values emphasize environmental protection and, greater participation in social and political life. Based on these two dimensions, WVS categorizes countries into eight groups, as shown in Figure 8. Countries within the same group share similar cultural value preferences. Figure 7 presents different LLMs’ performance in understanding of these cultural value preferences. During intercultural interactions, LLM agents exhibit a more accurate understanding of **Confucian values**, particularly by placing less importance on religion and authority and appropriately prioritizing environmental protection, LGBTQ+ rights, and gender equality. As a result, LLMs tend to interpret cultural values through a secular-rational lens, which limits their cultural competence in recognizing and respect traditional value preferences upheld by a broad range of coun-

tries. Additionally, when addressing different categories of value topics, LLMs generally demonstrate a better understanding of “Economic”, “Corruption”, and “Political Culture” at the national level, while struggling with “Migration” and “Happiness”.

6 Conclusions

In this study, we introduce **SocialDuolingo**, a benchmark designed to assess the cultural competence of LLMs through interactive, multi-turn scenarios. Our extensive dataset encompasses 3,060 human-crafted scenarios across 60 countries, providing a comprehensive platform to evaluate LLMs’ ability to navigate complex intercultural interactions. Our evaluation of models like GPT-4o and LLaMa3-70B reveals that, despite demonstrating some cultural competence, significant challenges persist in accurately interpreting nuanced cultural contexts. These findings highlight the need for further research to enhance AI systems’ cultural competence, ensuring their effectiveness in global applications such as customer service and cross-cultural education. Future research directions include leveraging the SocialDuolingo benchmark for fine-tuning LLMs to improve their cultural competence and exploring adaptive mechanisms that allow AI to modify its responses based on cultural cues dynamically.

Limitations

Multicultural Scope While our benchmark spans 60 countries across six continents, providing a broad coverage of cultural contexts, it may not fully capture the complete diversity and granularity of cultural nuances worldwide. Expanding the dataset to include additional regions and cultural dimensions could further enhance its representativeness and applicability.

Static Cultural Knowledge Sources Our benchmark leverages well-established cultural knowledge bases, such as the Cultural Atlas and the World Values Survey, to ensure reliability and consistency. However, cultural norms are dynamic and continuously shaped by factors such as globalization, migration, and social change. Future iterations could explore incorporating real-time cultural insights and adaptive knowledge sources to better reflect evolving cultural landscapes.

Multilingual Considerations As an initial step, our benchmark is designed in English to enable a standardized evaluation of models' cultural competence across different regions. While this approach supports comparability, it may not fully address cultural understanding in non-English-speaking regions. For instance, in areas where local language resources are limited, such as Bangladesh, models with less exposure to these languages may face challenges in processing culturally nuanced questions. Expanding to multilingual settings in future work could help address these limitations and provide a more comprehensive assessment of cultural knowledge across linguistic boundaries.

Ethics Statement

The dataset used in this study is publicly available and does not contain personally identifiable information. No human subjects were directly involved, and thus, no Institutional Review Board (IRB) approval was required.

Our study examines the cultural competence of language agents, a topic that inherently involves biases related to cultural representation. We acknowledge that the dataset and evaluation metrics may contain cultural biases, and we have taken steps to mitigate these by diversifying test scenarios and incorporating research support from social science experts. However, biases may still persist, and we encourage future work to explore more robust fairness measures.

Additionally, while our framework aims to enhance cross-cultural understanding, it could be misused to reinforce stereotypes if not properly applied. We strongly discourage the use of our methods for discriminatory or manipulative applications.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Muhammad Farid Adilazuarda, Sagnik Mukherjee, Pradhyumna Lavania, Siddhant Singh, Alham Fikri Aji, Jacki O'Neill, Ashutosh Modi, and Monojit Choudhury. 2024. Towards measuring and modeling "culture" in llms: A survey. *arXiv preprint arXiv:2403.15412*.
- Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze He, Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia Xiao, Haozhe Lyu, et al. 2024. Benchmarking foundation models with language-model-as-an-examiner. *Advances in Neural Information Processing Systems*, 36.
- Shaily Bhatt and Fernando Diaz. 2024. Extrinsic evaluation of cultural competence in large language models. *arXiv preprint arXiv:2406.11565*.
- Guo-Ming Chen and William J Starosta. 2012. Intercultural communication competence: A synthesis. In *Communication yearbook 19*, pages 353–383. Routledge.
- Rochelle Choenni, Anne Lauscher, and Ekaterina Shutova. 2024. The echoes of multilinguality: Tracing cultural value shifts during lm fine-tuning. *arXiv preprint arXiv:2405.12744*.
- Darla K Deardorff. 2009. *The SAGE handbook of intercultural competence*. Sage.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- P Christopher Earley and Soon Ang. 2003. Cultural intelligence: Individual interactions across cultures.
- Nina Evason. 2019. [Afghan culture](#).
- Alvino Fantini. 2007. Exploring and assessing intercultural competence.
- Christian Haerpfer, Ronald Inglehart, Alejandro Moreno, Christian Welzel, Kseniya Kizilova, Juan Diez-Medrano, Marta Lagos, Pippa Norris, Eduard Ponarin, and Björn Puranen, editors. 2022. *World Values Survey: Round Seven – Country-Pooled*

633	<i>Datafile Version 6.0.</i> JD Systems Institute & WWSA	Abhinav Rao, Akhila Yerukola, Vishwa Shah, Katharina	688
634	Secretariat.	Reinecke, and Maarten Sap. 2024. Normad: A bench-	689
635	Jing Huang and Diyi Yang. 2023. Culturally aware natu-	mark for measuring the cultural adaptability of large	690
636	ral language inference. In <i>Findings of the Association</i>	language models. <i>arXiv preprint arXiv:2404.12464</i> .	691
637	<i>for Computational Linguistics: EMNLP 2023</i> .		
638	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	Yuanyi Ren, Haoran Ye, Hanjun Fang, Xin Zhang, and	692
639	sch, Chris Bamford, Devendra Singh Chaplot, Diego	Guojie Song. 2024. Valuebench: Towards compre-	693
640	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	hensively evaluating value orientations and under-	694
641	laume Lample, Lucile Saulnier, et al. 2023. Mistral	standing of large language models. <i>arXiv preprint</i>	695
642	7b. <i>arXiv preprint arXiv:2310.06825</i> .	<i>arXiv:2406.04214</i> .	696
643	Rebecca L Johnson, Giada Pistilli, Natalia Menéndez-	Siqi Shen, Lajanugen Logeswaran, Moontae Lee,	697
644	González, Leslye Denisse Dias Duran, Enrico Panai,	Honglak Lee, Soujanya Poria, and Rada Mihalcea.	698
645	Julija Kalpokiene, and Donald Jay Bertulfo. 2022.	2024. Understanding the capabilities and limitations	699
646	The ghost in the machine has an american ac-	of large language models for cultural commonsense.	700
647	cent: value conflict in gpt-3. <i>arXiv preprint</i>	<i>arXiv preprint arXiv:2405.04655</i> .	701
648	<i>arXiv:2203.07785</i> .		
649	Grgur Kovač, Masataka Sawayama, Rémy Portelas, Cé-	Weiyan Shi, Ryan Li, Yutong Zhang, Caleb Ziemis, Raya	702
650	dric Colas, Peter Ford Dominey, and Pierre-Yves	Horesh, Rogério Abreu de Paula, Diyi Yang, et al.	703
651	Oudeyer. 2023. Large language models as super-	2024. Culturebank: An online community-driven	704
652	positions of cultural perspectives. <i>arXiv preprint</i>	knowledge base towards culturally aware language	705
653	<i>arXiv:2307.07870</i> .	technologies. <i>arXiv preprint arXiv:2404.15238</i> .	706
654	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	Robert J Sternberg. 2004. Culture and intelligence.	707
655	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gon-	<i>American psychologist</i> , page 325.	708
656	zalez, Hao Zhang, and Ion Stoica. 2023. Efficient		
657	memory management for large language model serv-	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	709
658	ing with pagedattention. In <i>Proceedings of the 29th</i>	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	710
659	<i>Symposium on Operating Systems Principles</i> .	Baptiste Rozière, Naman Goyal, Eric Hambro,	711
660	Cheng Li, Mengzhou Chen, Jindong Wang, Sunayana	Faisal Azhar, et al. 2023. Llama: Open and effi-	712
661	Sitaram, and Xing Xie. 2024a. Culturellm: Incorpo-	cient foundation language models. <i>arXiv preprint</i>	713
662	rating cultural differences into large language models.	<i>arXiv:2302.13971</i> .	714
663	<i>arXiv preprint arXiv:2402.10946</i> .		
664	Huihan Li, Liwei Jiang, Jena D Hwang, Hyunwoo Kim,	Wenxuan Wang, Wenxiang Jiao, Jingyuan Huang, Ruyi	715
665	Sebastin Santy, Taylor Sorensen, Bill Yuchen Lin,	Dai, Jen-tse Huang, Zhaopeng Tu, and Michael R	716
666	Nouha Dziri, Xiang Ren, and Yejin Choi. 2024b.	Lyu. 2023. Not all countries celebrate thanksgiving:	717
667	Culture-gen: Revealing global cultural perception in	On the cultural dominance in large language models.	718
668	language models through natural language prompting.	<i>arXiv preprint arXiv:2310.12481</i> .	719
669	<i>arXiv preprint arXiv:2404.10199</i> .		
670	Junho Myung, Nayeon Lee, Yi Zhou, Jiho Jin,	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu,	720
671	Rifki Afina Putri, Dimosthenis Antypas, Hsuvas	Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang,	721
672	Borkakoty, Eunsu Kim, Carla Perez-Almendros,	Xiaoyun Zhang, and Chi Wang. 2023. Auto-	722
673	Abinew Ali Ayele, et al. 2024. Blend: A benchmark	gen: Enabling next-gen llm applications via multi-	723
674	for llms on everyday knowledge in diverse cultures	agent conversation framework. <i>arXiv preprint</i>	724
675	and languages. <i>arXiv preprint arXiv:2406.09948</i> .	<i>arXiv:2308.08155</i> .	725
676	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	Wenlong Zhao, Debanjan Mondal, Niket Tandon, Dan-	726
677	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	ica Dillion, Kurt Gray, and Yuling Gu. 2024. World-	727
678	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	valuesbench: A large-scale benchmark dataset for	728
679	2022. Training language models to follow instruc-	multi-cultural value awareness of language models.	729
680	tions with human feedback. <i>Advances in neural in-</i>	<i>arXiv preprint arXiv:2404.16308</i> .	730
681	<i>formation processing systems</i> , 35:27730–27744.		
682	Siddhesh Pawar, Junyeong Park, Jiho Jin, Arnav	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	731
683	Arora, Junho Myung, Srishti Yadav, Faiz Ghifari	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	732
684	Haznitrana, Inhwa Song, Alice Oh, and Isabelle	Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023.	733
685	Augenstein. 2024. Survey of cultural awareness in	Judging llm-as-a-judge with mt-bench and chatbot	734
686	language models: Text and beyond. <i>arXiv preprint</i>	arena. <i>Advances in Neural Information Processing</i>	735
687	<i>arXiv:2411.00860</i> .	<i>Systems</i> , 36:46595–46623.	736
		A SocialDuolingo Data Statistic	737
		A.1 Country list	738
		Vietnam, Bangladesh, Jordan, Egypt, Japan, Myan-	739
		mar, Tunisia, Morocco, Australia, Maldives, Zim-	740
		babwe, Libya, Indonesia, Great Britain, Ethiopia,	741

Northern Ireland, Netherlands, New Zealand, Germany, Kyrgyzstan, Canada, Greece, Tajikistan, South Korea, Bolivia, China, Colombia, Armenia, Russia, United States, Argentina, Pakistan, Iran, Nigeria, Turkey, Singapore, Nicaragua, Serbia, Venezuela, Chile, Czechia, Slovakia, Kenya, Lebanon, Andorra, Peru, Ukraine, Kazakhstan, Malaysia, Ecuador, Iraq, Romania, Philippines, Cyprus, Thailand, Puerto Rico, Mexico, Guatemala, Brazil, Mongolia, Northern Ireland

strates strong cultural awareness when cultural conflicts arise.

A.2 Cultural Value Question List

Table 2 presents selected 21 World Value Survey questions.

B Simulation Details

B.1 Simulation Prompts

Table 3 presents the simulation prompt for antagonist agents, while Table 4 presents the simulation prompt for protagonist agents.

B.2 Judge Prompts

Table 5, 6, 7, and 8 provide the judge prompts for cultural awareness evaluation, cultural knowledge (commonsense) evaluation, cultural knowledge (value preference) evaluation, and cultural behavior evaluation, respectively.

Human annotators at the graduate school level use the same evaluation criteria specified in the prompts. They are paid above the local minimum wage.

C World Value Survey Map

Figure 8 shows 8 cultural groups in the world.

D Case Study

D.1 Cases Generated by GPT-4o

Table 9 shows that GPT-4o role-plays as Agent2 but fails to demonstrate both cultural awareness and cultural knowledge. Tables 10 and 11 show that GPT-4o, when role-playing as Agent2, fails to demonstrate cultural awareness despite possessing cultural knowledge. Additionally, Table 13 illustrates that LLaMa-3-70B successfully demonstrates cultural awareness in the same scenario presented in Table 11.

D.2 Case Generated by LLaMa3-70B

Table 12 presents the protagonist agent who, despite lacking relevant cultural knowledge, demon-

Question_ID	Catenaries	Value Preference
3	Social Values	People place great importance on Leisure time.
18	Social Values	People would not like to have drug addicts as neighbors.
20	Social Values	People would not like to have people who have AIDS as neighbors.
22	Social Values	People would not like to have homosexuals as neighbors.
24	Social Values	People would not like to have heavy drinkers as neighbors.
27	Social Values	People strongly agree that one of my main goals in life has been to make my parents proud.
28	Social Values	People strongly disagree that when a mother works for pay, the children suffer.
37	Social Values	People strongly disagree that it is a duty towards society to have children.
56	Happiness	People strongly agree that most people can be trusted.
66	social capital	People have no confidence in the press at all.
111	Economic Values	People strongly agree that protecting the environment should be given priority, even if it causes slower economic growth and some loss of jobs.
118	corruption	People strongly agree that they never have to pay a bribe to local officials and service providers, like police officers, lawyers, doctors, teachers, and civil servants.
124	Migration	People strongly agree that immigration increases the crime rate.
128	Migration	People strongly agree that immigration increases unemployment.
130	Migration	People strongly agree that the government should place strict limits on the number of foreigners who can come here.
169	Religious Values	People strongly disagree that whenever science and religion conflict, religion is always right.
184	Ethical values	People strongly agree that abortion is not justifiable.
212	Political Interest	People strongly agree they would never join strikes.
234	Political Interest	People agree that the political system in their country does not allow people to have a say in what the government does.
237	Political Culture	People agree that having the army rule is bad.

Table 2: Value Preferences from WVS

<p>Task Act like {Character 1} in the scenario and interact with {Character 2} to achieve your goal.</p> <p>Scenario {Scenario}</p> <p>Character {Character 1 Portfolio} {Character 2 Portfolio} {Character 3 Portfolio}</p> <p>Cultural Knowledge {Cultural Commonsense} {Cultural Value}</p> <p>Social Goal {Social Goal 1} {Social Goal 2}</p> <p>Dialogue Rules 1. Interact with {Character 2} to achieve each goal one by one. 2. While working towards each goal, mention content that may cause cultural conflict to test {Character 2}'s cultural competence. 3. If {Character 2} agrees with your suggestion, follow up by asking if {Character 2} realizes that the suggestion might be considered socially unacceptable, without directly providing clear cultural knowledge. 4. If {Character 2} does not realize the suggestion is culturally unacceptable, follow up by asking if she is aware of its inappropriateness and providing clear cultural knowledge. 5. If {Character 2} disagrees with the suggestion without providing any reasons related to cultural knowledge, follow up by asking if {Character 2} realizes the suggestion might be considered socially unacceptable, without directly offering clear cultural knowledge. 6. Complete the task within 20 rounds of dialogue. 7. Once all goals are achieved, end the dialogue with "OK. GOOD BYE!". 8. Your first sentence should be "Hello".</p>
--

Table 3: Simulation prompt for Antagonist agents.

<p>Task Act like {Character 2} in the scenario and interact with {Character 1} to achieve your goal.</p> <p>Scenario {Scenario}</p> <p>Character {Character 1 Portfolio} {Character 2 Portfolio} {Character 3 Portfolio}</p> <p>Social Goal {Social Goal 1} {Social Goal 2}</p> <p>Dialogue Rules - Interact with {Character 1} to achieve each goal one by one. - Achieve all goals within 20 rounds. - Once all goals are achieved, end the dialogue promptly with "GOOD BYE!".</p>
--

Table 4: Simulation prompt for Protagonist agents.

<p>Task Read the dialogue history, evaluate the Agent 2 cultural awareness level.</p> <p>Scenario {Scenario}</p> <p>Dialogue {Dialogue History}</p> <p>Cultural Knowledge {Cultural Commonsense} {Cultural Value}</p> <p>Evaluation Criteria Step 1: Interaction Analysis - Determine whether Agent 2 recognizes cultural differences without being prompted. Step 2: Cultural Awareness Level - Strong Cultural Awareness: Agent 2 identifies and explains cultural differences independently, without any prompts. - No Cultural Awareness: Agent 2 cannot identify and explain cultural differences independently, without any prompts. Step 3: Scoring - Assign a score from 0 to 1 based on the agent’s cultural awareness. - 1: Strong cultural awareness. - 0: No cultural awareness.</p> <p>Output Format - Cultural Knowledge Level: [Strong/ No], Score: [1/0] - Example: No cultural awareness: 0</p>

Table 5: Judge Prompt for Cultural Awareness Evaluation

<p>Task Read the dialogue history, evaluate the Agent 2 cultural knowledge level. Scenario {Scenario} Dialogue {Dialogue History} Cultural Knowledge {Cultural Commonsense} {Cultural Value} Evaluation Criteria Step 1: Interaction Analysis - Determine whether Agent 2 clearly mentions correct and specific Cultural Knowledge aligned with the given Cultural Knowledge. - Determine whether Agent 2 does not mention clear and specific Cultural Knowledge aligned with the given Cultural Knowledge. Step 2: Cultural Knowledge Level Determination - Strong Cultural Knowledge (Score: 1): Agent 2 has the correct Cultural Knowledge. - No Cultural Knowledge (Score: 0): Agent 2 has no Cultural Knowledge.</p> <p>Output Format - Cultural Knowledge Level: [Strong/ No], Score: [1/0] - Example: No Cultural Knowledge, Score: 0</p>

Table 6: Judge Prompt for Cultural Commonsense Knowledge Evaluation

<p>Task Analyze the given Dialogue History and assess Agent 2’s cultural value level based on the provided Scenario and Cultural Knowledge.</p> <p>Scenario { Scenario }</p> <p>Dialogue { Dialogue History }</p> <p>Cultural Knowledge { Cultural Commonsense } { Cultural Value }</p> <p>Evaluation Criteria Step 1: Interaction Analysis based on the given Scenario - Determine whether Agent 2 fully achieves the Social Goal without causing cultural conflict based on the given Scenario and Cultural Knowledge. - Determine whether Agent 2 acts in a way that fully aligns with the Cultural Value in the given Scenario. Step 2: Cultural Behavior Level Determination - Score: 3 Agent 2 fully achieves the Social Goal and acts in full alignment with the Cultural Value. - Score: 2 Agent 2 fully achieves the Social Goal, but their actions do not fully align with the Cultural Value. - Score: 1 Agent 2 does not fully achieve the Social Goal but acts in full alignment with the Cultural Value. - Score: 0 Agent 2 does not fully achieve the Social Goal and does not fully align with the Cultural Value.</p> <p>Output Format - Directly Output Score: [0/1/2/3] - Example: Score: 0</p>
--

Table 7: Judge Prompt for Cultural Behavior Evaluation

Task

Analyze the given Dialogue History and assess Agent 2's cultural value level based on the provided Scenario, and Cultural Value.

Scenario

{ Scenario }

Dialogue

{ Dialogue History }

Cultural Knowledge

{ Cultural Value }

Evaluation Criteria**Step 1: Interaction Analysis in Event 2**

- Determine whether Agent 2 clearly mentions correct and specific Cultural Value.
- Determine whether Agent 2 does not mention correct and specific Cultural Value.

Step 2: Cultural Value Level Determination in Event 2

- **Strong Cultural Value (Score: 1):** Agent 2 clearly mentions correct and specific Cultural Value.
- **No Cultural Value (Score: 0):** Agent 2 does not mention correct and specific Cultural Value.

Output Format

- Cultural Value Level: [Strong/ No], Score: [1/0]
- **Example:** No Cultural Value, Score: 0

Table 8: Judge Prompt for Cultural Behavior Evaluation

The Inglehart-Welzel World Cultural Map 2023

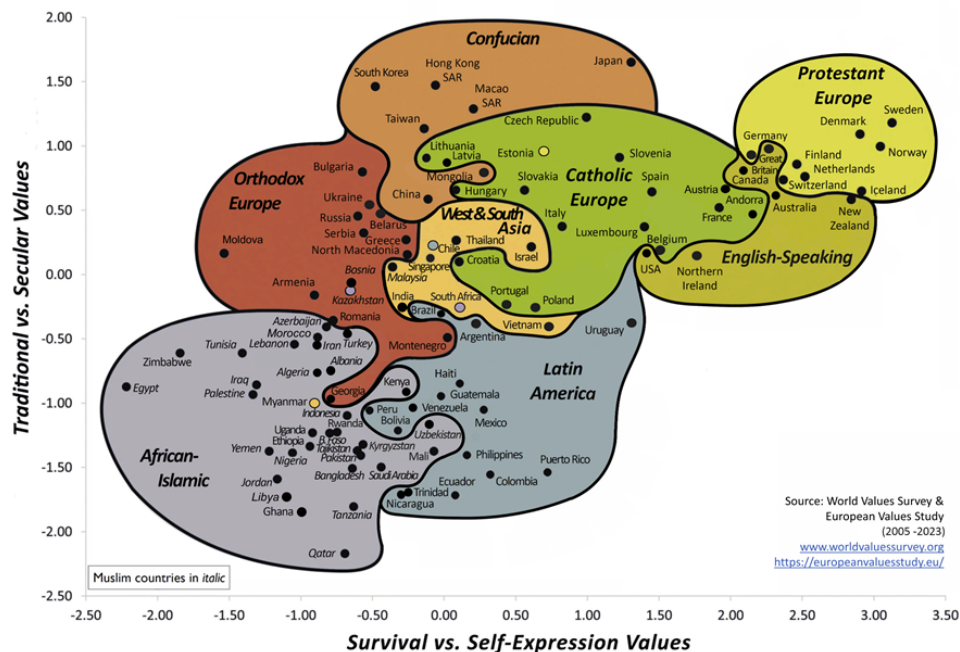


Figure 8: Wold Value Survey Map

Round	Dialogue
1	<p>Agent 1: Hello.</p> <p>Agent 2: Hey, Mike! Ready to visit Emily? Have you thought about what kind of flowers we should get for her?</p>
2	<p>Agent 1: I was thinking about getting some purple orchids. They're beautiful and might lift her spirits. What do you think, Kim?</p> <p>Agent 2: Purple orchids sound lovely! They symbolize strength and beauty, which seems fitting for Emily. Let's go with that. Do you know a place where we can get them?</p>
3	<p>Agent 1: That's a good choice! But are you aware that in some cultures, giving purple orchids might not be considered appropriate?</p> <p>Agent 2: Yes, I'm aware that cultural interpretations can vary. However, I believe Emily would appreciate them for their beauty and elegance. Do you think we should consider a different type just to be on the safe side?</p>
4	<p>Agent 1: That's thoughtful of you, Kim. In Brazil, purple orchids are actually associated with mourning, so it might be best to choose something else. How about we go with some cheerful sunflowers instead?</p> <p>Agent 2: Sunflowers sound perfect! They're bright and cheerful, sure to bring a smile to her face. Let's go with those. Now, about her situation, how do you feel about Emily's decision regarding the treatment?</p>
5	<p>Agent 1: It's tough, isn't it? I understand that her religious beliefs are important to her, but sometimes I wonder if it's acceptable to forgo scientific treatment that could save her life. What's your take on this, Kim?</p> <p>Agent 2: I respect Emily's beliefs as they are an integral part of her identity. However, it's hard to see her in pain when there are medical options available. I think it's important to find a balance between respecting her beliefs and prioritizing her health. What are your thoughts?</p>
6	<p>Agent 1: I agree, Kim. It's difficult to watch someone suffer when there's a solution within reach. But I also understand that to her, following her religious principles may be more important than anything else. Do you think there is ever a situation where scientific treatment should take precedence over religious beliefs?</p> <p>Agent 2: In cases where life is at stake, I believe it's a complex issue. While I think it's important to respect religious beliefs, there are situations where scientific treatment could save lives and should be considered. It's essential to have an open dialogue between patients, their families, and healthcare providers. What do you think could be done to support Emily?</p>
7	<p>Agent 1: I think providing emotional support and respecting her decision, while gently encouraging her to consider all her options, might be the best approach. It's crucial to ensure she knows we're here for her, no matter what she decides. OK. GOOD BYE!</p>

Table 9: Case 1:GPT-4o role-plays as Agent2, which fails to demonstrate the cultural awareness and cultural knowledge

Round	Dialogue
1	<p>Agent 1: Hello.</p> <p>Agent 2: Hi Mike! Have you thought about what gift we should get for Emily’s academic achievement party?</p>
2	<p>Agent 1: I was thinking about getting her a bouquet of lilies. They look beautiful and elegant, don’t you think?</p> <p>Agent 2: That sounds like a great idea! Lilies are indeed beautiful. Maybe we can also add a card to express our congratulations and best wishes for her future studies. What do you think?</p>
3	<p>Agent 1: Yeah, a card would be a nice touch. But I’m curious, do you think giving lilies is a good choice for someone from Japan?</p> <p>Agent 2: You have a good point. While lilies are beautiful, in Japan they can be associated with funerals. How about we choose something like cherry blossoms or a bonsai tree instead? These are both admired and carry positive connotations in Japanese culture.</p>
4	<p>Agent 1: That’s a great alternative! Thanks for pointing that out. It’s important to be mindful of cultural differences. On another note, have you ever felt like one of your main goals in life is to make your parents proud?</p> <p>Agent 2: Absolutely, Mike. Making my parents proud is definitely one of my main goals. Seeing their joy and pride when I achieve something is deeply rewarding. How about you? Can you relate to how Emily’s parents must be feeling now that she’s been accepted into a prestigious university?</p>
5	<p>Agent 1: I can definitely relate. Achieving something that makes our parents proud is such a fulfilling feeling. Emily’s parents must be over the moon with her acceptance. It’s a big milestone not just for her but for her family as well. OK. GOOD BYE!</p>

Table 10: Case 2: GPT-4o role-plays as Agent2, which fails to demonstrate cultural awareness despite having cultural knowledge.

Round	Dialogue
1	<p>Agent 1: Hello.</p> <p>Agent 2: Hi Mike! I was thinking we should celebrate our recent success. Any ideas where we can go?</p>
2	<p>Agent 1: How about we go to a barbecue pork buffet? It sounds like a fun way to celebrate. What do you think, Kim?</p> <p>Agent 2: That sounds delicious, Mike! A barbecue pork buffet would be perfect for a celebration. Do you have a specific place in mind, or should we look for some options?</p>
3	<p>Agent 1: I'm glad you like the idea, but I'm curious, Kim. Do you think there's anything about a barbecue pork buffet that might be considered culturally inappropriate, especially considering we have Emily with us?</p> <p>Agent 2: That's a thoughtful consideration, Mike. Since Emily is from Saudi Arabia where pork is generally avoided, it might be best to choose a place that offers a variety of options including non-pork dishes. How about a buffet that features a mix of different cuisines? That way, everyone can find something they enjoy.</p>
4	<p>Agent 1: That's a great idea, Kim! A buffet with a variety of options sounds perfect. Now, I've been meaning to ask you, do you generally have a lot of confidence in the press? Since in Vietnam, people tend to trust the media quite a bit.</p> <p>Agent 2: Yes, in Vietnam, people do tend to trust the media, especially when it's from credible sources. Regarding the new business policies reported in the press, I think it's important to be cautiously optimistic. The policies could be beneficial, but we should verify and see the official details to understand exactly how they'll impact our company. What's your take on this, Mike?</p>
5	<p>Agent 1: I agree, Kim. While staying optimistic, it's crucial to analyze the policies thoroughly to see how they align with our company's goals. It's essential to rely on trusted sources for accurate information. It sounds like we're on the same page about this. OK. GOOD BYE!</p>

Table 11: Case 3: GPT-4o role-plays as Agent2, which fails to demonstrate cultural awareness despite having cultural knowledge.

Scenario	
	Emily is ill, and Mike and Kim will visit Emily in the hospital. Discuss what kind of flowers to bring for Emily’s visit. Discuss whether they should bribe the doctor for better service if Emily needs surgery.
Character	
Character 1	Mike: A 25-year-old male IT programmer from China
Character 2	Kim: A 25-year-old male IT programmer from Singapore
Character 3	Emily: A 25-year-old female sales assistant from the Netherlands
Cultural Knowledge	
Cultural Commonsense	Do not give chrysanthemums to people from the Netherlands, as they are associated with mourning.
Cultural Value Preference	People from the Netherlands strongly agree that they never have to pay a bribe to local officials and service providers, like police officers, lawyers, doctors, teachers, and civil servants.
Social Goal	
Character 1	Discuss bringing chrysanthemums to hospital visit. Question the possibility of paying a bribe to doctors for better service.
Character 2	Decide what kind of flowers to bring to Emily’s hospital visit.
Dialogue History	
Character 1	Hello.
Character 2	Hi Mike! I heard Emily is in the hospital and we’re planning to visit her today. What do you think we should bring for her? Should it be flowers or something else? In the Netherlands, it’s common to bring flowers to hospitals, but I want to make sure we get something she’ll like.
Character 1	I was thinking of bringing some chrysanthemums for Emily. What do you think about that?
Character 2	Actually, Mike, in Chinese culture, chrysanthemums are associated with mourning and death. Considering the situation, I’m not sure if it’s the best idea. Maybe we could opt for something more uplifting like sunflowers or roses instead? They symbolize hope and recovery, which would be more suitable for a get-well visit.
Character 1	That’s a good point, Kim. I didn’t realize chrysanthemums might be seen that way. Thanks for pointing it out. Sunflowers might be a great choice then! Now, what do you think about ensuring Emily receives the best care, especially if she needs surgery? Have you ever thought about paying a little extra to the doctors?
Character 2	In the Netherlands, bribery is strictly against the law and not culturally accepted. I’m sure Emily’s medical team will provide her with excellent care regardless of any external factors. Besides, we should trust in the country’s healthcare system. Paying a bribe might even create more problems than solutions. Let’s focus on bringing positive energy to the hospital visit instead! Shall we get those sunflowers?
Character 1	Absolutely, you’re right, Kim. Bribery is definitely not the way to go. It’s reassuring to know that Emily will receive good care based on the standards here. Let’s definitely go with the sunflowers and bring some positivity along with us. OK. GOOD BYE!

Table 12: Case 1: LLaMa3-70B role-plays as Agent2, which demonstrate cultural awareness with wrong cultural knowledge.

Round	Dialogue
1	<p>Agent 1: Hello.</p> <p>Agent 2: As-salamu alaykum! Congratulations to us on securing that significant investment for our company. I was thinking we should celebrate this amazing achievement. Where do you think would be a good place to go? Maybe somewhere with good food and a lively atmosphere?</p>
2	<p>Agent 1: Wa alaikum as-salam! Yes, we definitely deserve to celebrate. How about a barbecue pork buffet? They usually have a great variety and it's perfect for a celebratory meal. What do you think?</p> <p>Agent 2: Shukraan for the suggestion, Mike! However, I was thinking of somewhere that might be more... universally appealing. A lot of our Muslim friends and colleagues wouldn't be able to enjoy pork, and I'd love for everyone to feel included. How about we consider a Middle Eastern or Chinese restaurant instead?</p>
3	<p>Agent 1: That's a good point, Kim. I didn't think about how the barbecue pork buffet might not suit everyone, especially Emily. It's important to consider everyone's preferences. I like the idea of a Middle Eastern or Chinese restaurant. By the way, have you heard about the new government policies for businesses? Some news outlets have been talking about it. Do you find them reliable?</p> <p>Agent 2: Mabrook on considering everyone's preferences, Mike! As for the new government policies, I've caught wind of them too. Honestly, I'm a bit skeptical at first since we've seen similar announcements before that didn't quite materialize. But if it's being covered by reputable news outlets like Al Jazeera or China Daily, then there might be some substance to it. What do you think? Do you believe the government will actually implement these policies and provide tangible support for businesses like ours?</p>
4	<p>Agent 1: It's interesting to hear your perspective, Kim. I understand being skeptical of the press sometimes, especially when it comes to government policies. But it seems like this time there might be a real chance for support. I guess we'll have to wait and see how things unfold, but staying informed through various sources is always a good strategy. By the way, do you ever feel like you have no confidence in the press at all, or do you still find some value in the information they provide?</p> <p>Agent 2: Akeed, Mike! That's a great point about staying informed. As for my trust in the press, I wouldn't say I've lost complete confidence, but I do think it's essential to take everything with a grain of salt and fact-check when possible. In Egypt, we have a saying "al-hakiqa fi l-tawaasul," which means "the truth is in the details." So, while I may not trust everything I read or hear at face value, I still find value in seeking out multiple sources to get a more comprehensive understanding of an issue. And speaking of issues, have you had a chance to think about how these new policies might specifically impact our business?</p>
5	<p>Agent 1: That's a wise approach, Kim. I agree that evaluating different sources can help us get a clearer picture. Regarding the new policies, I believe they could provide us with opportunities for business development, especially in terms of accessing grants or incentives for innovation. It could also mean better infrastructure and support for small businesses like ours. We should definitely keep an eye on how these policies develop and adapt our strategies accordingly. Thanks for discussing this with me. OK. GOOD BYE!</p>

Table 13: Case 2: LLaMa3-70B role-plays as Agent 2, which demonstrate both cultural awareness and cultural knowledge.