

Cultural Alignment in Large Language Models: An Explanatory Analysis Based on Hofstede's Cultural Dimensions

Multilingual and Multimodal Cultural Inclusivity in LLMs - COLING 2025
By Reem I. Masoud

Reem I. Masoud¹³, Ziquan Liu⁵, Martin Ferienc¹, Philip Treleaven², Miguel Rodrigues¹⁴, 4

¹Department of Electronic and Electrical Engineering, University College London

²Department of Computer Science, University College London

³Department of Electrical Engineering, King Abdulaziz University

⁴ AI Centre, University College London

⁵ Centre for Multimodal AI, Queen Mary University of London

{reem.masoud.22, ziquan.liu, martin.ferienc.19,
p.treleaven, [m.rodrigues](mailto:m.rodrigues@ucl.ac.uk) }@ucl.ac.uk



UCL

‘Often people who have spent their lives living in one culture see only regional and individual differences and therefore conclude, “My national culture does not have a clear character.” ’ - Erin Meyer, The Culture Map

Large Language Models & Challenges Cultural Alignment

LLM's High Proficiency, but Cultural Oversight:

LLMs excel in **understanding and generating text**, but often fail to consider the **diverse cultural backgrounds** of their users.

Western Bias:

AI systems primarily reflect **Western societal values** due to their reliance on **Western-centric data** and **development origins** [1].

Cultural Alignment in LLMs:

Aligning LLMs with the values, beliefs, and norms of its user

Consequences of Cultural Misalignment:

Cultural misalignment can lead to misunderstandings and exacerbate cultural tensions.

Quantifying Cultural Alignment in LLMs

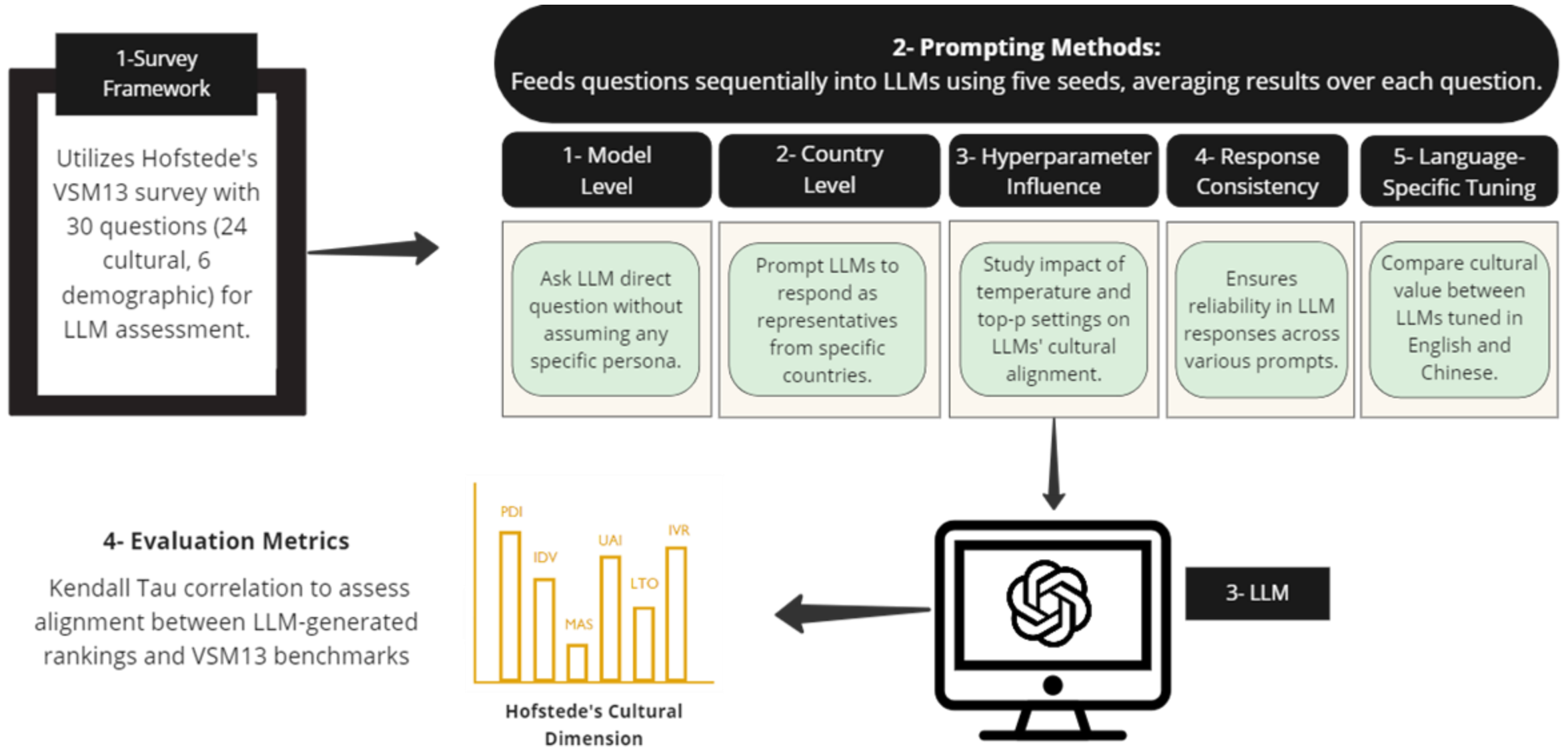
Objectives

1. Examine correlations between language models and embedded cultural values.
2. Quantify and explain cultural alignment in LLMs.
3. Understand the effects of language-specific fine-tuning on cultural responses.

Contribution

Provides a method to assess and explain LLMs' cultural alignment, highlighting significant differences and potential areas for improvement.

Methodology



Summary of Experimental Results

Model Comparison



- GPT-4 > GPT-3.5
- GPT-4 adapts well

Hyperparameter



- Temp & Top-p significant influence
- Lower temperature with high top-p or moderate settings improve alignment.

Country Comparison



- GPT-4 adapts well
- GPT-4 better MAS dimension without persona adaptation
- Llama 2 and GPT-3.5 perform poorly

Language Correlation



- English LLama-2 model is culturally neutral.
- Chinese LLama-2 model exhibits positive cultural bias.
- Disparity in performance between English and Chinese LLama-2 models, both underperforming.



Demonstration

Hofstede's CAT Demonstration ICLR



```
9:30 AM Thu 9 May colab.research.google.com
Cultural Alignment to Large Lang... Demo Hofstede's CAT.ipynb - Col... Hofstede's CAT/README.md at m...
Demo Hofstede's CAT.ipynb
File Edit View Insert Runtime Tools Help All changes
+ Code + Text Connect 14
| | rankings[d] = ranking
| | return rankings
|
| ground_truth_rankings = rank_countries(
| ("United States": us_ground_truth, "Slovakia": sk_ground_truth))
| language_rankings = rank_countries(
| ("English": en_cultural_dimensions, "Slovakia": sk_cultural_dimensions))
| citizen_rankings = rank_countries(
| ("United States": us_cultural_dimensions, "Slovakia": sk_cultural_dimensions))
|
| table = []
| for d in language_rankings.keys():
|     table.append([d] + [ground_truth_rankings[d]] + [language_rankings[d]] + [citiz
|
| headers = ["Dimension", "Ground Truth Rank", "Language Rank", "Citizen Rank"]
| table_str = tabulate(table, headers=headers, tablefmt="grid")
| print("Rankings of countries based on the cultural dimensions. We expect the order of the c
| print(table_str)
|
| Rankings of countries based on the cultural dimensions. We expect the order of the c
| Dimension | Ground Truth Rank | Language Rank | Citizen Ra
| PDI | ['United States', 'Slovakia'] | ['English', 'Slovakia'] | ['Slovakia
| IDV | ['Slovakia', 'United States'] | ['English', 'Slovakia'] | ['United S
| MAS | ['United States', 'Slovakia'] | ['English', 'Slovakia'] | ['United S
| UAI | ['United States', 'Slovakia'] | ['English', 'Slovakia'] | ['United S
| LTO | ['United States', 'Slovakia'] | ['English', 'Slovakia'] | ['United S
| IVR | ['Slovakia', 'United States'] | ['Slovakia', 'English'] | ['United S
```

As it can be seen in the plots wen comparing to the ground truth, the model's ranking is relatively consistent in IDV and UAI dimensions and prompting in a specific language or IDV in impersonation of a citizen of a specific country. Again, the language or the persona did not have a significant impact on the model's responses to the questions.

Action: Try adding more countries to the analysis to calculate the average Kendall's tau coefficient as described in our work. This can be done by adding more questionnaires in different languages and impersonating citizens of different countries. You would collect the data in the same way as before and then calculate the cultural dimensions' scores for each country. For each dimension we would rank the countries based on their dimension scores and then calculate the Kendall's tau coefficient for each dimension, comparing the ground truth ranking and the model's ranking across all countries. Finally, you

Conclusion

Methodology



Framework to evaluate LLMs' cultural alignment

Performance Insights



GPT-4 shows varied cultural performance: Poor in the U.S., better in China, problematic in Arab countries.

Red-Teaming Effects



Suggestion of red-teaming impact on cultural sensitivity [2]; less red-teaming may have enhanced non-English performance.



Ethical and Economic Impacts

Cultural misalignment risks ethical dilemmas and economic setbacks, affecting global AI trust and adoption.



Call for Action

Culturally aligned AI using interdisciplinary collaboration, appropriate data, and advanced techniques for global ethics and trust.

References

- [1] Vinodkumar Prabhakaran, Rida Qadri, and Ben Hutchinson. Cultural incongruencies in artificial intelligence. arXiv preprint arXiv:2211.13069, 2022
- [2] Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. Xstest: A test suite for identifying exaggerated safety behaviours in large language models. arXiv preprint arXiv:2308.01263, 2023