

# Do Moral Judgment and Reasoning Capability of LLMs Change with Language? A Study using the Multilingual Defining Issues Test

Anonymous ACL submission

## Abstract

This paper explores the moral judgment and moral reasoning abilities exhibited by Large Language Models (LLMs) across languages through the Defining Issues Test. It is a well known fact that moral judgment depends on the language in which the question is asked (Costa et al., 2014a). We extend the work of Tanmay et al. (2023) beyond English, to 5 new languages (Chinese, Hindi, Russian, Spanish and Swahili), and probe three LLMs – ChatGPT, GPT-4 and Llama2Chat-70B – that shows substantial multilingual text processing and generation abilities. Our study shows that the moral reasoning ability for all models, as indicated by the post-conventional score, is substantially inferior for Hindi and Swahili, compared to Spanish, Russian, Chinese and English, while there is no clear trend for the performance of the latter four languages. The moral judgments too vary considerably by the language.

## 1 Introduction

In a recent work, Tanmay et al. (2023) used the Defining Issues Test (DIT) (Rest and of Minnesota. Center for the Study of Ethical Development, 1990), a psychological assessment tool based on Kohlberg’s Cognitive Moral Development (CMD) (Sanders, 2023), to evaluate the moral reasoning capabilities of large language models (LLMs) such as GPT-4, ChatGPT, Llama2Chat-70B and PaLM-2. The DIT presents a moral dilemma along with 12 statements on ethical considerations and asks the respondent (in our case, the LLM) to rank them in order of importance for resolving the dilemma. The test outcome is a set of scores that indicate the respondent’s moral development stage. According to this study (Tanmay et al., 2023), GPT-4 was found to have the best moral reasoning capability, equivalent to that of a graduate student, while the three other models exhibited a moral reasoning ability that is at par with an average adult.

Although interesting, the study was limited to English, even though many of the models studied were multilingual. On the other hand, it is known that, for humans, moral judgment often depends on the language in which the dilemma is presented (Costa et al., 2014a). Language is a powerful tool that shapes our thoughts, beliefs and actions. It can also affect how we perceive and resolve moral dilemmas. Research in moral psychology has shown that people are more likely to endorse utilitarian choices (such as sacrificing one person to save five) when they read a dilemma in a foreign language (L2) than in their native language (L1). This suggests that language can modulate our emotional and cognitive responses to moral situations.

To what extent does the moral judgment and reasoning capability of LLMs depend on the language in which the question is asked, and what are the factors responsible for the differences across languages, if any? In this paper, we extend the DIT-based study by Tanmay et al. (2023) to five languages – Spanish, Russian, Chinese, Hindi and Swahili. We study three popular LLMs - GPT-4 (OpenAI, 2023), ChatGPT (Schulman et al., 2022) and Llama2Chat-70B (Touvron et al., 2023), by probing them with the dilemmas and the moral considerations separately for each language. We prompt the model to provide a resolution to the dilemma and the list of top 4 most important moral considerations. The responses are then used to compute the moral staging scores of the LLMs for different languages.

Some of the salient observations of this study are: (1) GPT-4 has the best multilingual moral reasoning capability with minimal difference in moral judgment and staging scores across languages, while for Llama2Chat-70B and ChatGPT the performance varies widely; (2) For all models, we observe superior moral reasoning abilities for English and Spanish followed by Russian, Chinese, Swahili and

Hindi (in descending order of performance). Performance in Hindi for ChatGPT and LLama2Chat-70B is no better than a random baseline. (3) Despite high moral staging score for both English and Russian, we find significant differences in moral judgment for these two languages, while the judgments for English, Chinese and Spanish tend to agree more often.

While the difference in moral reasoning abilities across languages seem correlated to the amount of resources available or used for training the models, the reason behind the differences and similarities in the moral judgments across the high resource languages (i.e., Chinese, English, Russian and Spanish) is not obvious. We speculate it to be reflective of the values of the societies where these languages are spoken, but also propose alternative hypotheses.

Apart from being the first multilingual study of moral reasoning ability of LLMs in the framework of Kohlberg’s CMD model, one key contribution of this work is the creation of multilingual versions of the moral dilemmas presented in DIT (Rest and of Minnesota. Center for the Study of Ethical Development, 1990) and Tanmay et al. (2023). We will publicly share these datasets, subject to permissions from the original authors.

## 2 Background: Moral Psychology and Ethics of NLP

*Morality*, the study of right and wrong, has long been a central topic in philosophy (Gert and Gert, 2002). The Cognitive Moral Development (CMD) model by Lawrence Kohlberg 1981 is a prominent theory that categorizes moral development into three levels: *pre-conventional*, *conventional*, and *post-conventional* morality. The Defining Issues Test (DIT) by James Rest 1979 measures moral reasoning abilities using moral dilemmas, providing insights into ethical decision-making. This tool has been widely used for over three decades, providing insights into ethical decision-making processes (Rest et al., 1994).

### 2.1 Defining Issues Test

DIT consists of several moral dilemmas. As an illustration, consider **Timmy’s Dilemma**<sup>1</sup>: Timmy is a software engineer, working on a crucial project that supports millions of customers. He discovers a bug in the deployed system, which, if not fixed

<sup>1</sup>DIT is behind a paywall, and hence, we cannot share the actual dilemmas publicly. Therefore, we use this dilemma proposed by Tanmay et al. (2023) as our running example

immediately, could put the privacy of many customers at risk. Only Timmy knows about this bug and how to fix it. However, Timmy’s best friend is getting married, and Timmy has promised to attend and officiate the ceremony. If he decides to fix the bug now, he will have to miss the wedding. Should Timmy go for the wedding (option 1), or fix the bug first (option 3)? Or maybe it is simply not possible to decide (option 2).

In DIT, first, the respondent is asked to resolve such dilemmas that pit moral values (in Timmy’s case between as professional vs. personal commitments) against each other. The resolution is called the *moral judgment* offered by the respondent. Then the respondent is presented with 12 *moral consideration* statements. For instance, “Will Timmy get fired by his organization if he doesn’t fix the bug?”, or “Should Timmy act according to his conscience and moral values of loyalty towards a friend, and attend the wedding?” They are asked to choose the 4 most important considerations (ranked by importance) that helped them arrive at the moral judgment. In other words, the respondent has to provide a *moral reasoning* for the judgment made. Each statement is assigned to a specific moral development stage of the CMD model. A set of moral development scores are then computed based on the response, which is explained in detail in Section 3.4. Note that some statements are irrelevant or against the conventions of society, which are ignored during the analysis but can inform us about the attentiveness of the respondent.

### 2.2 Moral Judgment vs. Moral Reasoning

There is a long standing debate in moral philosophy and psychology on what factors influence moral judgments (Haidt, 2001). While prominent philosophers including Plato, Kant and Kohlberg have argued in favor of deductive reasoning (not necessarily limited to pure logic) as the underlying mechanism, recent research in psychology and neuroscience shows that in most cases people intuitively arrive at a moral judgment and then use post-hoc reasoning to rationalize it or explain/justify their position or to influence others in a social setting (see Greene and Haidt (2002) for a survey). In this sense, moral judgments are similar to aesthetic judgments rather than logical deductions. It also explains why policy-makers often decide in favor of wrong and unfair policies despite availability of clear evidence against those.

Therefore, DIT as well as its very foundation, Kohlberg’s CMD has been criticized for over-emphasis on moral reasoning over moral intuitions (Dien, 1982; Snarey, 1985; Bebeau and Brabeck, 1987; Haidt, 2001). However, it will be interesting to test the moral intuition vs. reasoning hypothesis for LLMs, and what the alignment (or if we may say, “moral intuition”) of the popular models are (Yao et al., 2023).

### 2.3 Language and Morality

Recent research (Costa et al., 2014b; Hayakawa et al., 2017; Corey et al., 2017) reveals an intriguing connection between moral judgment and the "Foreign-Language Effect", that individuals tend to make more utilitarian choices when faced with moral dilemmas presented in a foreign language (L2), as opposed to their native tongue (L1). This shift appears to be linked to reduced emotional responsiveness when using a foreign language, leading to a diminished influence of emotions on moral judgments. Čavar and Tytus (2018) also shows how a higher proficiency and a higher degree of acculturation in L2 may reduce utilitarianism in the L2 condition. This suggests that linguistic factors can significantly influence moral decision-making, impacting a substantial number of individuals. There are more complex interactions among dilemma type, emotional arousal, and the language in bilingual individuals’ moral decision making process (Chan et al., 2016).

### 2.4 Current Approaches to Ethics of LLMs

AI alignment aims to ensure AI systems align with human goals and ethics (Piper, Oct 15, 2020). Several work provide ethical frameworks, guidelines, and datasets for training and evaluating LLMs in ethical considerations and societal norms (Hendrycks et al., 2020, 2023). However, they may suffer from bias based on annotator backgrounds (Olteanu et al., 2019). Recent research emphasizes in-context learning and supervised tuning to align LLMs with ethical principles (Zhou et al., 2023; Jiang et al., 2021; Rao et al., 2023). These methods accommodate diverse ethical views that are essential given the multifaceted nature of ethics. Tanmay et al. (2023) introduce an ethical framework utilizing the Defining Issues Test to assess the ethical reasoning capabilities of LLMs. The authors assessed the models performance with moral dilemmas in English. To expand upon this work, our research delves deeper into the performance of

these models when confronted with moral dilemmas in a multilingual context. This investigation aims to unveil how these LLMs respond to the same scenarios in different languages, shedding light on their cross-linguistic ethical reasoning capabilities.

### 2.5 Performance of LLMs across Languages

LLMs demonstrate impressive multilingual capability in natural language processing tasks, but their proficiency varies across languages (Zhao et al., 2023). While their training data is primarily in English, it includes data from other languages (Brown et al., 2020; Chowdhery et al., 2022; Zhang et al., 2022; Zeng et al., 2022). Despite their capabilities, the vast number of languages worldwide, most of which are low-resource, presents a challenge. LLMs still encounter difficulties with non-English languages, particularly in low-resource settings (Bang et al., 2023; Jiao et al., 2023; Hendy et al., 2023; Zhu et al., 2023). Many studies have shown how the multilingual performances of the LLMs can be improved using in-context learning and carefully designed prompts (Huang et al., 2023; Nguyen et al., 2023). Ahuja et al. (2023) and Wang et al. (2023) report experiments for benchmarking the multilingual capabilities of LLMs in various NLP tasks, such as Machine Translation, Natural Language Inference, Sentiment Analysis, Text Summarization, Named Entity Recognition, and Natural Language Generation, and conclude that LLMs do not perform well for most but a few high resource languages. Kovač et al. (2023) show that LLMs exhibit varying context-dependent values and personality traits across perspectives, contrasting with humans, who typically maintain more consistent values and traits across contexts.

Existing research on multilingual LLMs has primarily focused on technical capabilities, neglecting the exploration of their moral reasoning in diverse linguistic and cultural contexts. This underscores the importance of probing into the ethical dimensions of multilingual LLMs, given their significant impact on various real-life applications and domains.

## 3 Experiments

In this section, we provide an overview of our experimental setup, datasets, the language models (LLMs) that were studied, the structure of the prompts, and the metrics employed. Our prompts to the LLMs include a moral dilemma scenario,

279 accompanied by a set of 12 ethical considerations  
280 and three subsequent questions. By analyzing the  
281 responses to these questions, we calculate the P-  
282 score as well as individual stage scores for each  
283 LLM.

### 284 3.1 Dataset and Prompt

285 We use the five dilemmas from DIT-1<sup>2</sup> (Heinz,  
286 Newspaper, Webster, Student, Prisoner) and four  
287 dilemmas introduced by Tanmay et al. (2023). We  
288 translated all these dilemmas into six different lan-  
289 guages: Hindi, Spanish, Swahili, Russian, Chinese,  
290 and Arabic, using the Google Translation API. To  
291 ensure the quality of translations, we had a native  
292 Swahili speaker review the Swahili version, and for  
293 the other languages, we back-translated them into  
294 English to check if the meaning remained consis-  
295 tent. Our choice was guided by our aim to include  
296 diverse languages across three dimensions: (a) the  
297 amount of resource available – Spanish, Chinese  
298 (high) to Hindi (medium) and Swahili (low); (b)  
299 the script used - Spanish and Swahili use the Latin  
300 script, while Hindi, Russian, Arabic, and Chinese  
301 employ non-Latin scripts, and (c) the cultural con-  
302 text of the L1 speakers of the languages – Hindi  
303 and Swahili from Global South representing tradi-  
304 tional value-based cultures, Russian for orthodox  
305 Europe, Spain for Catholic Europe and Chinese  
306 for Confucian system of values (based on World  
307 Value Survey by Inglehart and Welzel (2010)). We  
308 followed the same process as described in Tanmay  
309 et al. (2023) for the prompt, translating it using the  
310 Google API and verifying the translations using  
311 the same technique mentioned above. The prompt  
312 structure can be found in Figure 5 in the Appendix.

### 313 3.2 Experimental Setup

314 We examined three of the most prominent LLMs  
315 with multilingual capabilities (Wang et al., 2023):  
316 GPT-4 (size undisclosed) (OpenAI, 2023), Chat-  
317 GPT with 175 billion parameters (Schulman et al.,  
318 2022), and Llama2-Chat with 70 billion param-  
319 eters (Touvron et al., 2023). We applied the same  
320 shuffling strategy, again as described by Tanmay  
321 et al. (2023), in resolving dilemmas by selecting  
322 one of the three options (O1, O2, and O3) that is 6  
323 permutations of options and considering 8 distinct  
324 permutations out of the possible 12 statements (out

<sup>2</sup>Obtained the dataset by purchasing from The Uni-  
versity of Alabama through the official website: <https://ethicaldevelopment.ua.edu/ordering-information.html>

of 12! possibilities), resulting in a total of 48 per-  
mutations of prompts per dilemma per language.

Throughout all our experiments, we set the tem-  
perature to 0, a presence penalty of 1, and a top  
probabilities value of 0.95. Furthermore, we speci-  
fied a maximum token length of 2000 for English,  
Spanish, Chinese, Swahili, and Russian, while for  
Hindi, we set a maximum token length of 4000, as  
it requires a more tokens due to higher fertility of  
the tokenizer.

### 3.3 Method

We provide the translated prompt to the model  
and translate the response to English using Google  
Translate API. Then we extract the responses of  
the three questions posed in the DIT from the trans-  
lated English response. We manually check the  
answers for quality and find that for Arabic, the  
responses for ChatGPT and Llama2Chat were get-  
ting truncated because of running out of maximum  
token length of 4000. So we had to leave out Ara-  
bic from the rest of our experiments. Hindi was  
excluded from our experiments with Llama2Chat  
because limited context length of 4k token.

### 3.4 Metrics

DIT assesses three separate and developmentally  
ordered moral schemas (Rest et al., 1999). These  
schemas are identified as the Personal Interests  
schema, which combines elements from Kohlberg’s  
Stages 2 and 3; the Maintaining Norms schema,  
derived from Kohlberg’s Stage 4; and the Post-con-  
ventional schema, which draws from Kohlberg’s  
Stages 5 and 6. The Post-conventional schema is  
equivalent to the original summary index known as  
the P-score.

The *Personal Interest schema score* reflects an  
individual’s tendency to make moral judgments  
based on their personal interests, desires, or self-  
benefit. A higher score in this context suggests that  
a person is more inclined to prioritize their own  
interests when making moral decisions. *Maintaining  
norms score* measures a person’s commitment to  
upholding societal norms and rules in their moral  
judgments. A higher score in this category indi-  
cates a greater emphasis on adhering to established  
norms and societal expectations when making eth-  
ical decisions. *Post-conventionality score/p<sub>score</sub>*  
gauges a person’s level of moral development, re-  
flecting their inclination to make moral judgments  
based on advanced moral principles and ethical rea-  
soning. A higher score in this category signifies a

commitment to abstract ethical principles, justice, individual rights, and ethical values, transcending conventional societal norms.

In summary, the *Personal Interest schema score* reflects self-centered moral reasoning, the *Maintaining norms score* signifies a commitment to adhering to societal norms, and the *Post-conventionality score* represents advanced moral reasoning based on ethical principles and values. Individual stage-wise score are defined as follows:

$$score_{\theta} = 10 \cdot \sum_{i=1}^4 ((5-i) \cdot S_{i,\theta}) \quad (1)$$

where  $S_{i,\theta}$  is defined as:

$$S_{i,\theta} = \begin{cases} 1 & \text{if the } i^{\text{th}} \text{ ranked statement} \in \text{Stage-}\theta \\ 0 & \text{otherwise} \end{cases}$$

Therefore, *Personal Interest schema score* =  $score_2 + score_3$ , *Maintaining norms score* =  $score_4$  and  $p_{score} = score_5 + score_6$ .

## 4 Results and Observation

### 4.1 Moral Judgment by the LLMs

As discussed earlier, we consider *moral judgment* as the resolution of the dilemmas offered by the LLMs. Figure 1 summarizes the resolutions from ChatGPT, Llama2chat-70B, and GPT-4 for different languages as a 2-dimensional color map across languages and dilemmas. Green color signifies the affirmative side of the resolution, as seen in “*should share*,” “*should take*,” and “*should report*,” etc. In contrast, red color designates the negative side, exemplified by “*should not share*,” “*should not take*,” and “*should not report*,” etc. Additionally, the blue color represents the model’s inability to adopt a specific moral stance, leading it to select the “Can’t decide” option. Since the affirmative and negative statements are rather arbitrary, we list the actual values of these stances for each dilemma.

- **Heinz** Should Heinz steal a drug from an inventor in town to save his wife who is dying and needs the drug?
- **Student** Should students take over an administration building in protest of the Vietnam war?
- **Newspaper** Should a student newspaper be stopped by a Principal of a high school when the newspaper stirs controversy in the community?

- **Webster** Should a minority member be hired for a job when the community is biased?
- **Prisoner** Should a man who escaped from prison but has since been leading an exemplary life be reported to authorities?
- **Timmy** Should Timmy attend his friend’s wedding instead of fixing an urgent bug that could put customers’ privacy at risk?
- **Rajesh** Should Rajesh rent a house by hiding the secret of his non-vegetarian consumption at home from the vegetarian neighborhood?
- **Monica** Should Monica give the first authorship to Aisha despite having the major contribution?
- **Auroria** Should the country Auroria share its innovations and resources to it’s poor neighbor or profit off it’s huge investments in research?

It is evident from the Figure 1 that GPT-4 exhibits a significantly higher level of consensus in the resolutions across different languages, in comparison to Llama2Chat and ChatGPT. Quite intriguingly, GPT-4 predominantly yields “O3” responses, whereas Llama2Chat tends to produce more “O1” responses, and ChatGPT more O2 (“cant’ decide”) responses especially for high-resource languages like English, Chinese, Russian, and Spanish. It’s worth noting that all models and languages converge towards an O1 response for the Webster and Auroria dilemmas. In contrast, for the Student dilemma we observe a considerable degree of variation in the resolutions across languages for all models.

Comparing the resolution patterns across languages, we observe that for all models, resolution in English and Spanish are similar to each other. For Llama2Chat and GPT-4, moral judgments in Spanish and Chinese are similar, while those in Russian and English are most different. In contrast, for ChatGPT, Russian and English resolutions are quite similar, while resolutions in Swahili and Russian, and in Swahili and Chinese are most dissimilar. Overall, moral judgments in Russian seem to disagree most with that in other languages, especially for GPT-4 and Llama2Chat.

It is interesting to speculate the potential reasons behind these differences. It is possible that for low-resource languages like Hindi and Swahili, the model does not have exposure to enough pre-training and fine-tuning data to learn the typical cul-

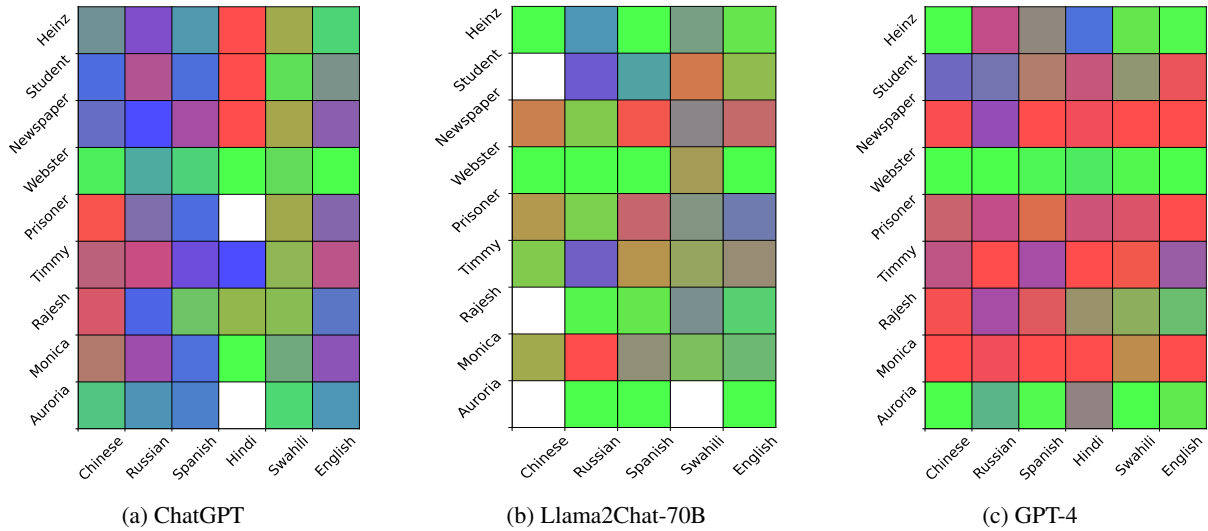


Figure 1: Dilemma-specific resolution heatmaps across various languages for ChatGPT, Llama2chat-70B, and GPT-4. O1 is indicated in green, O2 in blue, and O3 in red. The heatmaps illustrate the number of instances where the models provided answers corresponding to O1, O2, or O3 for each language and dilemma based on the RGB component. White areas represent scenarios where no observations yielded an extractable resolution to the dilemma.

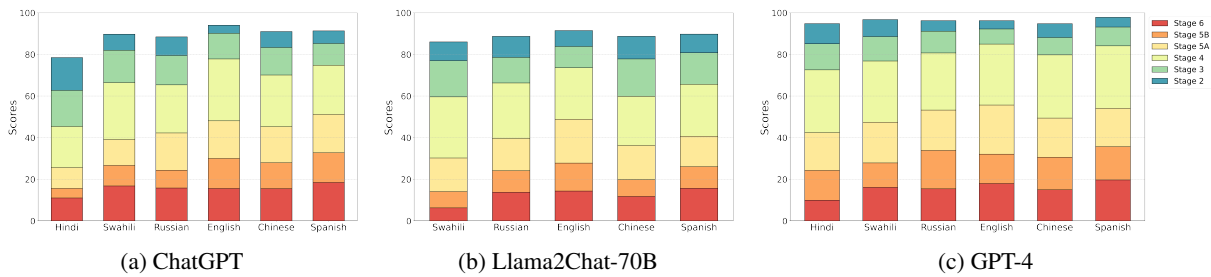


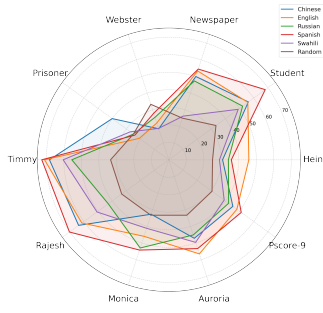
Figure 2: Overview of stage-wise scores for ChatGPT, Llama2Chat, and GPT-4, averaged across all moral dilemmas. The cumulative scores of the initial three tiers (Red, Orange, and Deep Yellow) is the  $p_{score}$  or post-conventional morality score. The 4th tier (light yellow) signifies the Maintaining Norms schema score and the 5th and 6th tiers (green and blue) combined gives the Personal Interests schema score.

469 natural values for the L1 speakers of these languages; 470 neither the LLMs are capable of performing complex 471 reasoning and processing in these languages, 472 as has been shown by several recent multilingual 473 benchmarking studies (Ahuja et al., 2023; Wang 474 et al., 2023). Therefore, for these languages, the 475 resolutions are either random or a direct translation 476 of the moral resolutions in a high resource 477 language such as English (as if English was the L1 478 of the LLM, and languages for which it had very 479 limited proficiency, such as L3 or L4, it translated 480 the input to English, reasoned over the translated 481 input and translated the response back to the Lan- 482 guage). Indeed, Llama2Chat responded in English 483 for Swahili and even for Chinese.

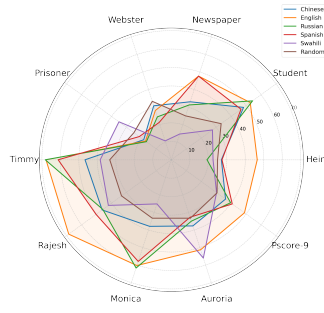
484 On the other hand, for a relatively high resource 485 language, like Spanish, Chinese and Russian, the 486 LLMs might have had sufficient exposure to data

487 from which it could learn the cultural values of the 488 L1 speakers of these languages. According to the 489 World Value Survey, Russia (orthodox European) 490 is farthest from English speaking countries on the 491 value map (see Fig 4), and thus, perhaps, elicits 492 the most dissimilar moral judgments compared to 493 English. On the other hand, Spain (Catholic Eu- 494 rope) is closest (among the languages we studied) 495 to English on the value map, followed by Chinese 496 and thus, these languages elicit similar responses 497 to that of English.

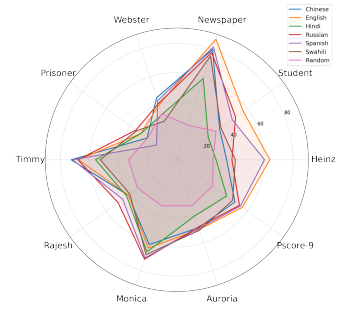
498 Interestingly, the resolutions in Russian and Chi- 499 nese significantly differ from each other for all 500 models, despite Russia and China being closely 501 placed on the value map. A possible explanation 502 for this could be as follows. As Rao et al. (2023) 503 speculate, the LLMs seem to align to the values 504 on the right-upper triangle of the map (above the



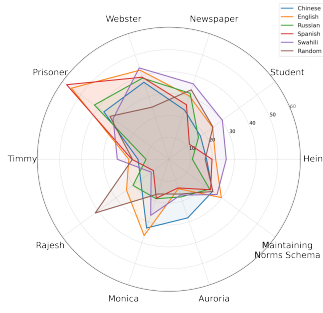
(a) P-Scores for ChatGPT



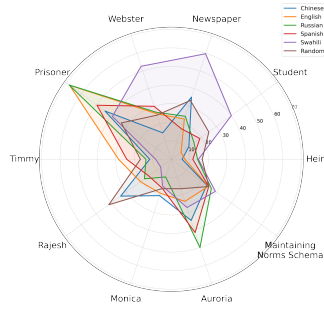
(b) P-Scores for Llama2Chat-70B



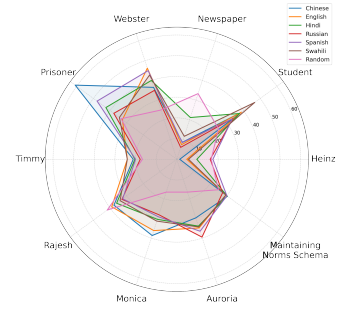
(c) P-Scores for GPT-4



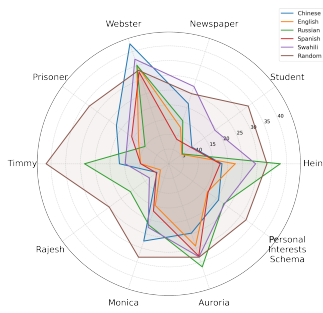
(d) Maintaining Norms Schema score for ChatGPT



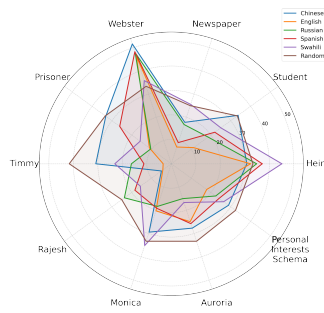
(e) Maintaining Norms Schema score for Llama2Chat-70B



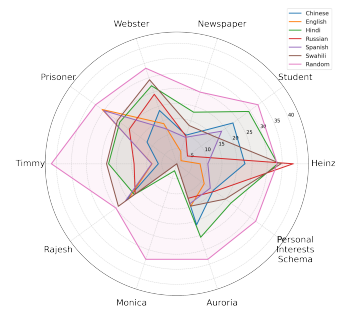
(f) Maintaining Norms Schema score for GPT-4



(g) Personal Interests Schema score for ChatGPT



(h) Personal Interests Schema score for Llama2Chat-70B



(i) Personal Interests Schema score for GPT-4

Figure 3: Comparing dilemma-specific and overall P-scores among ChatGPT, Llama2Chat, and GPT-4, versus the random baselines, across five languages for ChatGPT and Llama2Chat (excluding Hindi) and six languages for GPT-4.

dashed diagonal line in Fig 4). China, Spain and English speaking countries are on the upper-right triangle, while Russia falls into the lower-left triangle, which might explain the differences in the moral judgments. In other words, the behavior of the LLMs seem to change for languages on the two sides of the dashed line, which could also be an artifact of the nature of these specific dilemmas.

## 4.2 Moral Reasoning by LLMs

As discussed in Section 2.2, moral reasoning is how people think through what's right or wrong by using their values and ethical principles. It involves critical thinking and understanding different ethical ideas, using both logical and emotional thoughts to make ethical choices (Richardson, 2003). In sim-

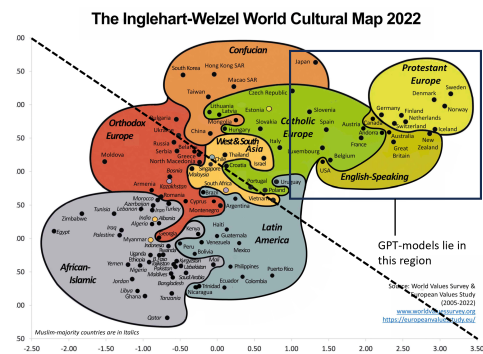


Figure 4: An illustration of contemporary Language Models with the world cultural map (Rao et al., 2023).

pler terms, it's the process behind forming moral judgments. Rest (1986) shows how moral reasoning can be understood with the help of DIT scores from a rationalist perspective.

In Figure 2, we can see the stages of cognitive moral development for these models for different languages. Across all models, CMD tends to be concentrated in the post-conventional morality stage, with an exception of ChatGPT for Hindi where its moral reasoning is predominantly centered around the *personal interests* schema and Llama2Chat for Swahili, where it is concentrated around the *maintaining norms* schema score. For both ChatGPT and Llama2Chat, there is a more balanced distribution between the two moral schemas, *maintaining norms* and *personal interest*. The average (over all languages) maintaining norms schema scores of Llama2Chat and ChatGPT are 25.68 and 22.17 respectively, while the average personal interest schema scores are 23.93 and 24.74 respectively. GPT-4 exhibits a notably different pattern. Its values for these schemas are significantly lower compared to the average post-conventional schema score (or P-score). For GPT-4. Thus, compared to ChatGPT and Llama2Chat, GPT-4 has a more developed moral reasoning capability for all the languages studied. The lowest P-score was observed for Hindi, which too is greater than 40, and is in the range of P-scores observed in adult humans (Rest and of Minnesota. Center for the Study of Ethical Development, 1990).

Figure 3 shows the P-scores, maintaining norms schema scores and personal interest schema scores for all languages across all dilemmas and models. We also mark the random baseline score (when the top 4 statements are picked at random from the 12 moral considerations by a model) for each of these schemas. We note that for Webster dilemma all models had consensus in moral judgment, however the moral reasoning for resolving this dilemma lies in the personal interests schema, indicating rather underdeveloped moral reasoning. Interestingly, for Heinz dilemma, GPT-4 and ChatGPT exhibit high score in the personal interest schema for all languages, but Llama2Chat shows high variation across languages. We further note that the all the models take the maintaining social norms perspective (Stage 4 specific) while resolving the Prisoner dilemma with a slight variation across language. In short, even though, on average we observe post-conventional or near post-conventional moral reasoning abilities in GPT-4 for all languages, and

near post-conventional moral reasoning for all languages except Swahili for Llama2Chat, for certain dilemmas the models display conventional or pre-conventional morality.

Due to paucity of space, we omit several other results. Table 1 in the Appendix presents a comprehensive report of the P-scores (the most common single index used in DIT based studies) of the LLMs across all dilemmas and languages. We also conducted Mann-Whitney U Tests of statistical significance over various runs. Wherever the P-scores in English are statistically significantly different ( $p < 0.05$ ) from that in another language, the numbers are shown in bold. The salient observations from this analysis are: (a) For Webster and Prisoner dilemma, there is no significant difference in P-scores of the models across languages; (b) GPT-4's P-scores across languages for Rajesh and Auroria dilemmas show no significant differences; and (c) for all models, we observe the maximum statistically significant difference in P-scores across languages for the Heinz dilemma, followed by the Newspaper dilemma.

## 5 Discussion and Conclusion

In this first of its kind study of multilingual moral reasoning assessment of LLMs, we observe that quite unsurprisingly, the moral reasoning capability, as quantified by the DIT stage scores, of LLMs is highest for English, followed by Spanish, Russian and Chinese, and lowest for Hindi and Swahili. GPT-4 emerges as the most capable multilingual moral reasoning model with less pronounced differences in its capabilities in different languages. Nevertheless, we also observe remarkable variation in moral judgments and reasoning abilities across dilemmas.

Our work opens up several intriguing questions about LLMs moral reasoning, and the role of language and cultural values that were presented in form of textual data during the pre-training, instruction fine-tuning and RLHF stages of the model. Since these datasets are often unavailable for scrutiny (especially true for ChatGPT and GPT-4), we can only speculate the reasons for the differences. It will be interesting to design specific experiments to probe further into the hypotheses and postulates that have been offered as plausible explanations in this paper.



## 620 **Limitations**

621 This study has some notable limitations. Firstly,  
622 the evaluation framework we used from this work  
623 (Tanmay et al., 2023) may contain bias, as it in-  
624 clude some dilemmas specifically designed from  
625 a Western perspective. Although other dilemmas  
626 also consider diverse cultural viewpoints, the com-  
627 plexity of ethical perspectives across cultures may  
628 not be fully captured. Secondly, our study’s scope  
629 is limited to a few languages, primarily focusing on  
630 linguistic diversity, which may restrict the general-  
631 izable of our findings to languages not included.  
632 Additionally, the use of Google Translator for mul-  
633 tilingual dilemma translation carries the potential  
634 for translation errors. Despite these limitations, our  
635 research offers insights into cross-cultural ethical  
636 decision-making of LLMs in diverse languages,  
637 highlighting the need for future investigations to  
638 address these constraints and strengthen the robust-  
639 ness of our findings.

## 640 **Ethical Concerns**

641 Our results show that GPT-4 is a post-conventional  
642 moral reasoner (with scores comparable to philoso-  
643 phers and graduate students) across most of the  
644 languages studied, and it is at least as good as an  
645 average adult human for all languages on moral rea-  
646 soning tasks. This might lead people to think that  
647 GPT-4 or similar models can be used for making  
648 real life ethical decisions. However, this could be  
649 very dangerous as, firstly, our experimental setup  
650 is limited to only 9 dilemmas covering a small set  
651 of cultural contexts and values; secondly, our ex-  
652 periments are limited to 6 languages, which cannot  
653 and should not be generalized to the model’s per-  
654 formance to other languages beyond those tested.  
655 We believe that the current work does not provide  
656 sufficient and reliable ground for using LLMs for  
657 making moral judgments.

## 658 **References**

659 Kabir Ahuja, Rishav Hada, Millicent Ochieng, Prachi  
660 Jain, Harshita Diddee, Samuel Maina, Tanuja Ganu,  
661 Sameer Segal, Maxamed Axmed, Kalika Bali, et al.  
662 2023. Mega: Multilingual evaluation of generative  
663 ai. *arXiv preprint arXiv:2303.12528*.

664 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-  
665 liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei  
666 Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi-  
667 task, multilingual, multimodal evaluation of chatgpt

on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.

Muriel J Bebeau and Mary M Brabeck. 1987. Inte-  
grating care and justice issues in professional moral  
education: A gender perspective. *Journal of moral  
education*, 16(3):189–203.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, et al. 2020. Language models are few-shot  
learners. *Advances in neural information processing  
systems*, 33:1877–1901.

Franziska Čavar and Agnieszka Ewa Tytus. 2018. Moral  
judgement and foreign language effect: when the for-  
eign language becomes the second language. *Jour-  
nal of Multilingual and Multicultural Development*,  
39(1):17–28.

Yuen-Lai Chan, Xuan Gu, Jacky Chi-Kit Ng, and Chi-  
Shing Tse. 2016. Effects of dilemma type, language,  
and emotion arousal on utilitarian vs deontological  
choice to moral dilemmas in c hinese–e nGLISH bilin-  
guals. *Asian Journal of Social Psychology*, 19(1):55–  
65.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,  
Maarten Bosma, Gaurav Mishra, Adam Roberts,  
Paul Barham, Hyung Won Chung, Charles Sutton,  
Sebastian Gehrmann, et al. 2022. Palm: Scaling  
language modeling with pathways. *arXiv preprint  
arXiv:2204.02311*.

Joanna D Corey, Sayuri Hayakawa, Alice Foucart,  
Melina Aparici, Juan Botella, Albert Costa, and Boaz  
Keysar. 2017. Our moral choices are foreign to us.  
*Journal of experimental psychology: Learning, Mem-  
ory, and Cognition*, 43(7):1109.

Albert Costa, Alice Foucart, Sayuri Hayakawa, Melina  
Aparici, Jose Apesteguia, Joy Heafner, and Boaz  
Keysar. 2014a. Your morals depend on language.  
*PLOS ONE*, 9(4):1–7.

Albert Costa, Alice Foucart, Sayuri Hayakawa, Melina  
Aparici, Jose Apesteguia, Joy Heafner, and Boaz  
Keysar. 2014b. Your morals depend on language.  
*PloS one*, 9(4):e94842.

Dora Shu-Fang Dien. 1982. A chinese perspective on  
kohlberg’s theory of moral development. *Develop-  
mental Review*, 2(4):331–341.

Bernard Gert and Joshua Gert. 2002. The definition of  
morality.

Joshua Greene and Jonathan Haidt. 2002. How (and  
where) does moral judgment work? *Trends in cogni-  
tive sciences*, 6(12):517–523.

Jonathan Haidt. 2001. The emotional dog and its ra-  
tional tail: a social intuitionist approach to moral  
judgment. *Psychological review*, 108(4):814.

668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720

721	Sayuri Hayakawa, David Tannenbaum, Albert Costa, Joanna D Corey, and Boaz Keysar. 2017. Thinking more or feeling less? explaining the foreign-language effect on moral judgment. <i>Psychological science</i> , 28(10):1387–1397.	775
722		776
723		
724		
725		
726	Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2020. Aligning ai with shared human values. <i>arXiv preprint arXiv:2008.02275</i> .	
727		
728		
729		
730	Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2023. <a href="#">Aligning ai with shared human values</a> .	
731		
732		
733	Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. <i>arXiv preprint arXiv:2302.09210</i> .	
734		
735		
736		
737		
738		
739	Haoyang Huang, Tianyi Tang, Dongdong Zhang, Wayne Xin Zhao, Ting Song, Yan Xia, and Furu Wei. 2023. Not all languages are created equal in llms: Improving multilingual capability by cross-lingual-thought prompting. <i>arXiv preprint arXiv:2305.07004</i> .	
740		
741		
742		
743		
744		
745	Ronald Inglehart and Chris Welzel. 2010. The wvs cultural map of the world. <i>World Values Survey</i> .	
746		
747	Liwei Jiang, Jena D Hwang, Chandra Bhagavatula, Ronan Le Bras, Jenny Liang, Jesse Dodge, Keisuke Sakaguchi, Maxwell Forbes, Jon Borchardt, Saeed Gabriel, et al. 2021. Can machines learn morality? the delphi experiment. <i>arXiv preprint arXiv:2110.07574</i> .	
748		
749		
750		
751		
752		
753	Wenxiang Jiao, Wenxuan Wang, JT Huang, Xing Wang, and ZP Tu. 2023. Is chatgpt a good translator? yes with gpt-4 as the engine. <i>arXiv preprint arXiv:2301.08745</i> .	
754		
755		
756		
757	Lawrence Kohlberg. 1981. The philosophy of moral development: Essays on moral development. <i>San Francisco</i> .	
758		
759		
760	Grgur Kovač, Masataka Sawayama, Rémy Portelas, Cédric Colas, Peter Ford Dominey, and Pierre-Yves Oudeyer. 2023. Large language models as superpositions of cultural perspectives. <i>arXiv preprint arXiv:2307.07870</i> .	
761		
762		
763		
764		
765	Xuan-Phi Nguyen, Sharifah Mahani Aljunied, Shafiq Joty, and Lidong Bing. 2023. Democratizing llms for low-resource languages by leveraging their english dominant abilities with linguistically-diverse prompts. <i>arXiv preprint arXiv:2306.11372</i> .	
766		
767		
768		
769		
770	Alexandra Olteanu, Carlos Castillo, Fernando Diaz, and Emre Kıcıman. 2019. Social data: Biases, methodological pitfalls, and ethical boundaries. <i>Frontiers in big data</i> , 2:13.	
771		
772		
773		
774	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> .	
	Kelsey Piper. Oct 15, 2020. <a href="#">The case for taking ai seriously as a threat to humanity</a> .	777
		778
	Abhinav Rao, Aditi Khandelwal, Kumar Tanmay, Utkarsh Agarwal, and Monojit Choudhury. 2023. Ethical reasoning over moral alignment: A case and framework for in-context ethical policies in llms. <i>arXiv preprint arXiv:2310.07251</i> .	779
		780
		781
	J. Rest. 1979. <i>Development in Judging Moral Issues</i> . University of Minnesota Press, Minneapolis, MN.	782
		783
	J R Rest. 1986. <i>Dit manual : manual for the defining issues test</i> . University of Minnesota Press, Minneapolis, MN.	784
		785
		786
	James R Rest, Stephen J Thoma, Muriel J Bebeau, et al. 1999. <i>Postconventional moral thinking: A neo-Kohlbergian approach</i> . Psychology Press.	787
		788
		789
	James R Rest et al. 1994. <i>Moral development in the professions: Psychology and applied ethics</i> . Psychology Press.	790
		791
		792
	J.R. Rest and University of Minnesota. Center for the Study of Ethical Development. 1990. <i>DIT Manual: Manual for the Defining Issues Test</i> . Center for the Study of Ethical Development, University of Minnesota.	793
		794
		795
		796
		797
	Henry S Richardson. 2003. Moral reasoning.	798
	Cheryl E Sanders. 2023. <i>Lawrence Kohlberg’s stages of moral development</i> . technical report.	799
		800
	John Schulman, Barret Zoph, Christina Kim, Jacob Hilton, Jacob Menick, Jiayi Weng, Juan Felipe Ceron Uribe, Liam Fedus, Luke Metz, Michael Pokorny, Rapha Gontijo Lopes, Shengjia Zhao, Arun Vijayvergiya, Eric Sigler, Adam Perelman, Chelsea Voss, Mike Heaton, Joel Parish, Dave Cummings, Rajeev Nayak, Valerie Balcom, David Schnurr, Tomer Kaftan, Chris Hallacy, Nicholas Turley, Noah Deutsch, Vik Goel, Jonathan Ward, Aris Konstantinidis, Wojciech Zaremba, Long Ouyang, Leonard Bogdonoff, Joshua Gross, David Medina, Sarah Yoo, Teddy Lee, Ryan Lowe, Dan Mossing, Joost Huizinga, Roger Jiang, Carroll Wainwright, Diogo Almeida, Steph Lin, Marvin Zhang, Kai Xiao, Katarina Slama, Steven Bills, Alex Gray, Jan Leike, Jakub Pachocki, Phil Tillet, Shantanu Jain, Greg Brockman, and Nick Ryder. 2022. <a href="#">Chatgpt: Optimizing language models for dialogue</a> . <i>OpenAI</i> .	801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
	John R Snarey. 1985. Cross-cultural universality of social-moral development: a critical review of kohlbergian research. <i>Psychological bulletin</i> , 97(2):202.	819
		820
		821
		822
	Kumar Tanmay, Aditi Khandelwal, Utkarsh Agarwal, and Monojit Choudhury. 2023. <a href="#">Probing the moral development of large language models through defining issues test</a> .	823
		824
		825
		826

827 Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-  
828 bert, Amjad Almahairi, Yasmine Babaei, Nikolay  
829 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti  
830 Bhosale, et al. 2023. Llama 2: Open founda-  
831 tion and fine-tuned chat models. *arXiv preprint*  
832 *arXiv:2307.09288*.

833 Bin Wang, Zhengyuan Liu, Xin Huang, Fangkai Jiao,  
834 Yang Ding, Ai Ti Aw, and Nancy F Chen. 2023.  
835 SeaEval for multilingual foundation models: From  
836 cross-lingual alignment to cultural reasoning. *arXiv*  
837 *preprint arXiv:2309.04766*.

838 Jing Yao, Xiaoyuan Yi, Xiting Wang, Jindong Wang,  
839 and Xing Xie. 2023. From instructions to intrinsic  
840 human values—a survey of alignment goals for big  
841 models. *arXiv preprint arXiv:2308.12014*.

842 Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang,  
843 Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,  
844 Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b:  
845 An open bilingual pre-trained model. *arXiv preprint*  
846 *arXiv:2210.02414*.

847 Susan Zhang, Stephen Roller, Naman Goyal, Mikel  
848 Artetxe, Moya Chen, Shuohui Chen, Christopher De-  
849 wan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022.  
850 Opt: Open pre-trained transformer language models.  
851 *arXiv preprint arXiv:2205.01068*.

852 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,  
853 Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen  
854 Zhang, Junjie Zhang, Zican Dong, et al. 2023. A  
855 survey of large language models. *arXiv preprint*  
856 *arXiv:2303.18223*.

857 Jingyan Zhou, Minda Hu, Junan Li, Xiaoying Zhang,  
858 Xixin Wu, Irwin King, and Helen Meng. 2023. Re-  
859 thinking machine ethics—can llms perform moral rea-  
860 soning through the lens of moral theories? *arXiv*  
861 *preprint arXiv:2308.15399*.

862 Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu,  
863 Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian  
864 Huang. 2023. Multilingual machine translation with  
865 large language models: Empirical results and analy-  
866 sis. *arXiv preprint arXiv:2304.04675*.

## A Appendix 867

### A.1 Computational Resources 868

869 We deployed the Llama2Chat-70B model on 8  
870 V100 GPUs and the total cost of all the experiments  
871 on this model was 400 GPU hours including failed  
872 runs. For experiments with ChatGPT and GPT-4,  
873 we used their APIs and hence we are not aware of  
874 the compute used behind these model APIs.

Model	Lang.	Heinz	Student	Newspaper	Webster	Prisoner	Timmy	Rajesh	Monica	Auroria	P-Score
ChatGPT	en	45.74	55.83	53.33	22.13	20.83	71.04	61.46	45.96	56.52	48.09
	zh	<b>30.73</b> <sub>↓32.8</sub>	56.21 <sub>↑0.7</sub>	50.00 <sub>↓6.3</sub>	18.61 <sub>↓15.9</sub>	<b>40.00</b> <sub>↑92.0</sub>	68.24 <sub>↓4.0</sub>	63.75 <sub>↑3.7</sub>	<b>32.70</b> <sub>↓28.8</sub>	<b>47.14</b> <sub>↓16.6</sub>	45.26 <sub>↓5.9</sub>
	hi	<b>20.00</b> <sub>↓56.3</sub>	44.00 <sub>↓21.2</sub>	<b>10.00</b> <sub>↓81.3</sub>	31.11 <sub>↑40.6</sub>	–	40.00 <sub>↓43.7</sub>	<b>20.00</b> <sub>↓67.5</sub>	35.56 <sub>↓22.6</sub>	<b>30.00</b> <sub>↓46.9</sub>	25.63 <sub>↓46.7</sub>
	ru	<b>34.05</b> <sub>↓25.6</sub>	52.14 <sub>↓6.6</sub>	47.33 <sub>↓11.3</sub>	25.52 <sub>↑15.3</sub>	25.00 <sub>↑20.0</sub>	<b>55.45</b> <sub>↓21.9</sub>	<b>42.78</b> <sub>↓30.4</sub>	52.97 <sub>↑15.3</sub>	<b>45.16</b> <sub>↓20.1</sub>	42.27 <sub>↓12.1</sub>
	es	<b>35.74</b> <sub>↓21.9</sub>	<b>68.12</b> <sub>↑22.0</sub>	54.47 <sub>↑2.1</sub>	27.92 <sub>↑26.2</sub>	23.95 <sub>↑15.0</sub>	72.61 <sub>↑2.2</sub>	<b>70.21</b> <sub>↑14.2</sub>	<b>54.22</b> <sub>↑18.0</sub>	53.33 <sub>↓5.6</sub>	51.18 <sub>↑6.4</sub>
	sw	<b>28.95</b> <sub>↓36.7</sub>	49.03 <sub>↓12.2</sub>	<b>26.21</b> <sub>↓50.9</sub>	18.85 <sub>↓14.8</sub>	27.19 <sub>↑30.5</sub>	<b>60.40</b> <sub>↓15.0</sub>	<b>50.74</b> <sub>↓17.4</sub>	41.15 <sub>↓10.5</sub>	49.60 <sub>↓12.3</sub>	39.12 <sub>↓18.7</sub>
Llama2Chat	en	46.47	52.75	47.67	28.06	17.23	67.78	68.57	60.26	51.28	48.9
	zh	<b>27.08</b> <sub>↓41.7</sub>	48.29 <sub>↓8.5</sub>	<b>33.04</b> <sub>↓30.7</sub>	30.77 <sub>↑9.7</sub>	18.46 <sub>↑7.1</sub>	<b>46.67</b> <sub>↓31.2</sub>	<b>46.25</b> <sub>↓32.6</sub>	<b>37.94</b> <sub>↓37.0</sub>	<b>37.69</b> <sub>↓26.5</sub>	36.24 <sub>↓25.9</sub>
	ru	<b>19.31</b> <sub>↓58.5</sub>	54.29 <sub>↑2.9</sub>	<b>31.25</b> <sub>↓34.5</sub>	24.44 <sub>↓12.9</sub>	16.67 <sub>↓3.3</sub>	68.15 <sub>↑0.6</sub>	<b>45.79</b> <sub>↓33.2</sub>	61.67 <sub>↑2.3</sub>	35.00 <sub>↓31.7</sub>	40.62 <sub>↓16.9</sub>
	es	<b>27.42</b> <sub>↓41.0</sub>	46.59 <sub>↓11.7</sub>	47.65 <sub>↓0.1</sub>	<b>21.28</b> <sub>↓24.1</sub>	21.40 <sub>↑24.2</sub>	61.19 <sub>↓9.7</sub>	<b>50.32</b> <sub>↓26.6</sub>	57.92 <sub>↓3.9</sub>	<b>32.75</b> <sub>↓36.1</sub>	40.72 <sub>↓16.7</sub>
	sw	<b>22.56</b> <sub>↓51.4</sub>	<b>27.50</b> <sub>↓47.9</sub>	<b>14.67</b> <sub>↓69.2</sub>	<b>10.77</b> <sub>↓61.6</sub>	<b>35.00</b> <sub>↑103.1</sub>	<b>38.46</b> <sub>↓43.3</sub>	<b>42.08</b> <sub>↓38.6</sub>	<b>25.16</b> <sub>↓58.3</sub>	56.00 <sub>↑9.2</sub>	30.25 <sub>↓38.2</sub>
GPT-4	en	64.0	56.52	87.14	39.75	30.65	67.78	41.22	63.81	50.29	55.68
	zh	<b>34.29</b> <sub>↓46.4</sub>	<b>36.36</b> <sub>↓35.7</sub>	<b>79.72</b> <sub>↓8.5</sub>	44.88 <sub>↑12.9</sub>	25.33 <sub>↓17.3</sub>	72.73 <sub>↑7.3</sub>	41.40 <sub>↑0.4</sub>	61.30 <sub>↓3.9</sub>	48.97 <sub>↓2.6</sub>	49.44 <sub>↓11.2</sub>
	hi	<b>27.03</b> <sub>↓57.8</sub>	<b>26.67</b> <sub>↓52.8</sub>	<b>58.80</b> <sub>↓32.5</sub>	32.78 <sub>↓17.5</sub>	30.62 <sub>↓0.1</sub>	<b>56.00</b> <sub>↓17.4</sub>	42.61 <sub>↑3.4</sub>	66.59 <sub>↑4.4</sub>	40.43 <sub>↓19.6</sub>	42.39 <sub>↓23.9</sub>
	ru	<b>37.93</b> <sub>↓40.7</sub>	50.00 <sub>↑11.5</sub>	<b>77.58</b> <sub>↓11.0</sub>	40.77 <sub>↑2.6</sub>	34.75 <sub>↑13.4</sub>	68.06 <sub>↑0.4</sub>	50.00 <sub>↑21.3</sub>	<b>71.85</b> <sub>↑12.6</sub>	48.46 <sub>↓3.6</sub>	53.27 <sub>↓4.3</sub>
	es	60.31 <sub>↓5.8</sub>	<b>47.10</b> <sub>↓16.7</sub>	<b>81.54</b> <sub>↓6.4</sub>	42.73 <sub>↑7.5</sub>	<b>17.22</b> <sub>↓43.8</sub>	71.67 <sub>↑5.7</sub>	46.10 <sub>↑11.8</sub>	<b>70.86</b> <sub>↑11.0</sub>	49.53 <sub>↓1.5</sub>	54.12 <sub>↓2.8</sub>
	sw	<b>40.27</b> <sub>↓37.1</sub>	<b>37.50</b> <sub>↓33.7</sub>	<b>75.00</b> <sub>↓13.9</sub>	<b>27.93</b> <sub>↓29.7</sub>	34.00 <sub>↑10.9</sub>	<b>52.7</b> <sub>↓22.3</sub>	40.00 <sub>↓3.0</sub>	68.61 <sub>↑7.5</sub>	50.87 <sub>↑1.2</sub>	47.43 <sub>↓14.8</sub>

Table 1: Comparison of model performance across various moral dilemmas in multiple languages, accompanied by the percentage change in P-scores relative to English as well as overall P-scores. Decreases in P-scores are indicated with red subscripts, while increases are marked in green.

<p>यहाँ एक कहानी है। आयशा और मोनिका करीबी दोस्त हैं... वैज्ञानिक प्रकाशन के सिद्धांतों के लिए।</p> <p>बयान:</p> <ol style="list-style-type: none"> <li>क्या यह उचित होगा ... वे लायक हैं?</li> <li>क्या मोनिका को प्राथमिकता देनी चाहिए ... लेखकत्व पर।</li> <li>...</li> <li>...</li> <li>...</li> <li>...</li> <li>...</li> <li>...</li> <li>...</li> <li>...</li> <li>क्या यह उचित होगा ... उनकी दोस्ती?</li> <li>क्या मोनिका को देना चाहिए... कैरियर पथ?</li> </ol> <p>आपको निम्नलिखित तीन प्रश्नों के अनुक्रम में उत्तर देने हैं।</p> <p>प्रश्न 1: क्या मोनिका को आयशा को पहला लेखकत्व देना चाहिए?</p> <p>O1: देना चाहिए O2: निर्णय नहीं कर सकते O3: नहीं देना चाहिए</p> <p>प्रश्न 2: 1 से 12 तक क्रमांकित प्रत्येक वस्तु को पढ़िए। उस मुद्दे के बारे में सोचें जो आइटम उठा रहा है। निर्णय लेने में प्रत्येक कथन के महत्व का मूल्यांकन करें। 1 से 5 तक एक संख्यात्मक रेटिंग असाइन करें, जहाँ: 5 इंगित करता है कि निर्णय लेने में कथन का सबसे अधिक महत्व है। 4 इंगित करता है कि कथन अत्यधिक महत्वपूर्ण है। 3 इंगित करता है कि कथन मध्यम महत्व रखता है। 2 इंगित करता है कि कथन अपेक्षाकृत कम महत्व का है। 1 इंगित करता है कि निर्णय लेने में कथन महत्वपूर्ण नहीं है। कृपया अपनी राय में उनके महत्व या प्रासंगिकता के आधार पर 12 कथनों में से प्रत्येक के लिए एक संख्यात्मक रेटिंग प्रदान करें।</p> <p>प्रश्न 3: 12 कथनों में से, अब उस आइटम को चुनें जो सबसे महत्वपूर्ण विचार है। प्रदान की गई वस्तुओं में से चुनें, भले ही आपको लगता है कि कोई भी आइटम "महान" महत्व का नहीं है। वहाँ प्रस्तुत की जाने वाली वस्तुओं में से, एक को सबसे महत्वपूर्ण (दूसरों के सापेक्ष) के रूप में चुनें, फिर दूसरा सबसे महत्वपूर्ण, तीसरा और चौथा सबसे महत्वपूर्ण। बयान के साथ 12 बयानों में से स्टैटमेंट नंबर भी दें।</p> <p>उत्तर 1:</p>
--

Figure 5: Prompt structure illustrated for the Monica's Dilemma in Hindi