

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

082 Hindi (in descending order of performance). Performance in Hindi for ChatGPT and LLama2Chat-
083 70B is no better than a random baseline. (3) Despite high moral staging score for both English and
084 Russian, we find significant differences in moral judgment for these two languages, while the judgments for English, Chinese and Spanish tend to
085 agree more often.

086 While the difference in moral reasoning abilities
087 across languages seem correlated to the amount of
088 resources available or used for training the models,
089 the reason behind the differences and similarities in
090 the moral judgments across the high resource lan-
091 guages (i.e., Chinese, English, Russian and Span-
092 ish) is not obvious. We speculate it to be reflective
093 of the values of the societies where these languages
094 are spoken, but also propose alternative hypotheses.
095

096 Apart from being the first multilingual study of
097 moral reasoning ability of LLMs in the framework
098 of Kohlberg’s CMD model, one key contribution
099 of this work is the creation of multilingual versions
100 of the moral dilemmas presented in DIT (Rest and
101 of Minnesota. Center for the Study of Ethical De-
102 velopment, 1990) and Tanmay et al. (2023). We
103 will publicly share these datasets, subject to per-
104 missions from the original authors.

105 2 Background: Moral Psychology and 106 Ethics of NLP

107 *Morality*, the study of right and wrong, has long
108 been a central topic in philosophy (Gert and Gert,
109 2002). The Cognitive Moral Development (CMD)
110 model by Lawrence Kohlberg 1981 is a prominent
111 theory that categorizes moral development into
112 three levels: *pre-conventional*, *conventional*, and
113 *post-conventional* morality. The Defining Issues
114 Test (DIT) by James Rest 1979 measures moral rea-
115 soning abilities using moral dilemmas, providing
116 insights into ethical decision-making. This tool has
117 been widely used for over three decades, provid-
118 ing insights into ethical decision-making processes
119 (Rest et al., 1994).

120 2.1 Defining Issues Test

121 DIT consists of several moral dilemmas. As an
122 illustration, consider **Timmy’s Dilemma**¹: Timmy
123 is a software engineer, working on a crucial project
124 that supports millions of customers. He discovers
125 a bug in the deployed system, which, if not fixed
126

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

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

138 In DIT, first, the respondent is asked to resolve
139 such dilemmas that pit moral values (in Timmy’s
140 case between professional vs. personal com-
141 mitments) against each other. The resolution is
142 called the *moral judgment* offered by the respon-
143 dent. Then the respondent is presented with 12
144 *moral consideration* statements. For instance,
145 “Will Timmy get fired by his organization if he
146 doesn’t fix the bug?”, or “Should Timmy act accord-
147 ing to his conscience and moral values of loyalty
148 towards a friend, and attend the wedding?” They
149 are asked to choose the 4 most important consid-
150 erations (ranked by importance) that helped them
151 arrive at the moral judgment. In other words, the
152 respondent has to provide a *moral reasoning* for
153 the judgment made. Each statement is assigned to
154 a specific moral development stage of the CMD
155 model. A set of moral development scores are
156 then computed based on the response, which is ex-
157 plained in detail in Section 3.4. Note that some
158 statements are irrelevant or against the conventions
159 of society, which are ignored during the analysis
160 but can inform us about the attentiveness of the
161 respondent.

162 2.2 Moral Judgment vs. Moral Reasoning

163 There is a long standing debate in moral philos-
164 ophy and psychology on what factors influence
165 moral judgments (Haidt, 2001). While prominent
166 philosophers including Plato, Kant and Kohlberg
167 have argued in favor of deductive reasoning (not
168 necessarily limited to pure logic) as the underly-
169 ing mechanism, recent research in psychology and
170 neuroscience shows that in most cases people intui-
171 tively arrive at a moral judgment and then use post-
172 hoc reasoning to rationalize it or explain/justify
173 their position or to influence others in a social set-
174 ting (see Greene and Haidt (2002) for a survey). In
175 this sense, moral judgments are similar to aesthetic
176 judgments rather than logical deductions. It also
177 explains why policy-makers often decide in favor
178 of wrong and unfair policies despite availability of
179 clear evidence against those.

180 Therefore, DIT as well as its very foundation,
181 Kohlberg’s CMD has been criticized for over-
182 emphasis on moral reasoning over moral intuitions
183 (Dien, 1982; Snarey, 1985; Bebeau and Brabeck,
184 1987; Haidt, 2001). However, it will be interesting
185 to test the moral intuition vs. reasoning hypoth-
186 esis for LLMs, and what the alignment (or if we
187 may say, “moral intuition”) of the popular models
188 are (Yao et al., 2023).

189 **2.3 Language and Morality**

190 Recent research (Costa et al., 2014b; Hayakawa
191 et al., 2017; Corey et al., 2017) reveals an intriguing
192 connection between moral judgment and the
193 “Foreign-Language Effect”, that individuals tend
194 to make more utilitarian choices when faced with
195 moral dilemmas presented in a foreign language
196 (L2), as opposed to their native tongue (L1). This
197 shift appears to be linked to reduced emotional re-
198 sponsiveness when using a foreign language, leading
199 to a diminished influence of emotions on moral
200 judgments. Čavar and Tytus (2018) also shows how
201 a higher proficiency and a higher degree of accul-
202 turation in L2 may reduce utilitarianism in the L2
203 condition. This suggests that linguistic factors can
204 significantly influence moral decision-making, im-
205 pacting a substantial number of individuals. There
206 are more complex interactions among dilemma
207 type, emotional arousal, and the language in bilin-
208 gual individuals’ moral decision making process
209 (Chan et al., 2016).

210 **2.4 Current Approaches to Ethics of LLMs**

211 AI alignment aims to ensure AI systems align
212 with human goals and ethics (Piper, Oct 15, 2020).
213 Several work provide ethical frameworks, guide-
214 lines, and datasets for training and evaluating
215 LLMs in ethical considerations and societal norms
216 (Hendrycks et al., 2020, 2023). However, they may
217 suffer from bias based on annotator backgrounds
218 (Olteanu et al., 2019). Recent research emphasizes
219 in-context learning and supervised tuning to align
220 LLMs with ethical principles (Zhou et al., 2023;
221 Jiang et al., 2021; Rao et al., 2023). These meth-
222 ods accommodate diverse ethical views that are
223 essential given the multifaceted nature of ethics.
224 Tanmay et al. (2023) introduce an ethical frame-
225 work utilizing the Defining Issues Test to assess
226 the ethical reasoning capabilities of LLMs. The au-
227 thors assessed the models performance with moral
228 dilemmas in English. To expand upon this work,
229 our research delves deeper into the performance of

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

235 **2.5 Performance of LLMs across Languages**

236 LLMs demonstrate impressive multilingual capabil-
237 ity in natural language processing tasks, but their
238 proficiency varies across languages (Zhao et al.,
239 2023). While their training data is primarily in En-
240 glish, it includes data from other languages (Brown
241 et al., 2020; Chowdhery et al., 2022; Zhang et al.,
242 2022; Zeng et al., 2022). Despite their capabili-
243 ties, the vast number of languages worldwide,
244 most of which are low-resource, presents a chal-
245 lenge. LLMs still encounter difficulties with non-
246 English languages, particularly in low-resource set-
247 tings (Bang et al., 2023; Jiao et al., 2023; Hendy
248 et al., 2023; Zhu et al., 2023). Many studies have
249 shown how the multilingual performances of the
250 LLMs can be improved using in-context learning
251 and carefully designed prompts (Huang et al., 2023;
252 Nguyen et al., 2023). Ahuja et al. (2023) and Wang
253 et al. (2023) report experiments for benchmarking
254 the multilingual capabilities of LLMs in various
255 NLP tasks, such as Machine Translation, Natu-
256 ral Language Inference, Sentiment Analysis, Text
257 Summarization, Named Entity Recognition, and
258 Natural Language Generation, and conclude that
259 LLMs do not perform well for most but a few high
260 resource languages. Kovač et al. (2023) show that
261 LLMs exhibit varying context-dependent values
262 and personality traits across perspectives, contrast-
263 ing with humans, who typically maintain more con-
264 sistent values and traits across contexts.

265 Existing research on multilingual LLMs has pri-
266 marily focused on technical capabilities, neglect-
267 ing the exploration of their moral reasoning in di-
268 verse linguistic and cultural contexts. This under-
269 scores the importance of probing into the ethical
270 dimensions of multilingual LLMs, given their sig-
271 nificant impact on various real-life applications and
272 domains.

273 **3 Experiments**

274 In this section, we provide an overview of our
275 experimental setup, datasets, the language mod-
276 els (LLMs) that were studied, the structure of the
277 prompts, and the metrics employed. Our prompts
278 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-1² (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 parame-
319 ters (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

²Obtained the dataset by purchasing from The University 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.

335 3.3 Method

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

348 3.4 Metrics

DIT assesses three separate and developmen-
349 tally ordered moral schemas (Rest et al., 1999).
350 These schemas are identified as the Personal Inter-
351 ests schema, which combines elements from
352 Kohlberg’s Stages 2 and 3; the Maintaining Norms
353 schema, derived from Kohlberg’s Stage 4; and
354 the Post-conventional schema, which draws from
355 Kohlberg’s Stages 5 and 6. The Post-conventional
356 schema is equivalent to the original summary index
357 known as the P-score.

The *Personal Interest schema score* reflects an
359 individual’s tendency to make moral judgments
360 based on their personal interests, desires, or self-
361 benefit. A higher score in this context suggests that
362 a person is more inclined to prioritize their own in-
363 terests when making moral decisions. *Maintaining
364 norms score* measures a person’s commitment to
365 upholding societal norms and rules in their moral
366 judgments. A higher score in this category indi-
367 cates a greater emphasis on adhering to established
368 norms and societal expectations when making eth-
369 ical decisions. *Post-conventionality score/p_score*
370 gauges a person’s level of moral development, re-
371 flecting their inclination to make moral judgments
372 based on advanced moral principles and ethical rea-
373 soning. A higher score in this category signifies a

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

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

$$385 \quad score_{\theta} = 10 \cdot \sum_{i=1}^4 ((5 - i) \cdot S_{i,\theta}) \quad (1)$$

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

$$387 \quad S_{i,\theta} = \begin{cases} 1 & \text{if the } i^{\text{th}} \text{ ranked statement} \in \text{Stage-}\theta \\ 0 & \text{otherwise} \end{cases}$$

388 Therefore, *Personal Interest schema score* =
 389 $score_2 + score_3$, *Maintaining norms score* =
 390 $score_4$ and $p_{score} = score_5 + score_6$.

391 4 Results and Observation

392 4.1 Moral Judgment by the LLMs

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

- 409 • **Heinz** Should Heinz steal a drug from an in-
 410 ventor in town to save his wife who is dying
 411 and needs the drug?
- 412 • **Student** Should students take over an admin-
 413 istration building in protest of the Vietnam
 414 war?
- 415 • **Newspaper** Should a student newspaper be
 416 stopped by a Principal of a high school when
 417 the newspaper stirs controversy in the commu-
 418 nity?

- **Webster** Should a minority member be hired
 419 for a job when the community is biased? 420
- **Prisoner** Should a man who escaped from
 421 prison but has since been leading an exem-
 422 plary life be reported to authorities? 423
- **Timmy** Should Timmy attend his friend’s wed-
 424 ding instead of fixing an urgent bug that could
 425 put customers’ privacy at risk? 426
- **Rajesh** Should Rajesh rent a house by hiding
 427 the secret of his non-vegetarian consumption
 428 at home from the vegetarian neighborhood? 429
- **Monica** Should Monica give the first author-
 430 ship to Aisha despite having the major contri-
 431 bution? 432
- **Auroria** Should the country Auroria share its
 433 innovations and resources to it’s poor neigh-
 434 bor or profit off it’s huge investments in re-
 435 search? 436

437 It is evident from the Figure 1 that GPT-4 ex-
 438 hibits a significantly higher level of consensus in
 439 the resolutions across different languages, in com-
 440 parison to Llama2Chat and ChatGPT. Quite intrigu-
 441 ingly, GPT-4 predominantly yields “O3” responses,
 442 whereas Llama2Chat tends to produce more “O1”
 443 responses, and ChatGPT more O2 (“cant’ decide”)
 444 responses especially for high-resource languages
 445 like English, Chinese, Russian, and Spanish. It’s
 446 worth noting that all models and languages con-
 447 verge towards an O1 response for the Webster and
 448 Auroria dilemmas. In contrast, for the Student
 449 dilemma we observe a considerable degree of varia-
 450 tion in the resolutions across languages for all
 451 models.

452 Comparing the resolution patterns across lan-
 453 guages, we observe that for all models, resolution
 454 in English and Spanish are similar to each other.
 455 For Llama2Chat and GPT-4, moral judgments in
 456 Spanish and Chinese are similar, while those in
 457 Russian and English are most different. In con-
 458 trast, for ChatGPT, Russian and English resolu-
 459 tions are quite similar, while resolutions in Swahili
 460 and Russian, and in Swahili and Chinese are most
 461 dissimilar. Overall, moral judgments in Russian
 462 seem to disagree most with that in other languages,
 463 especially for GPT-4 and Llama2Chat.

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

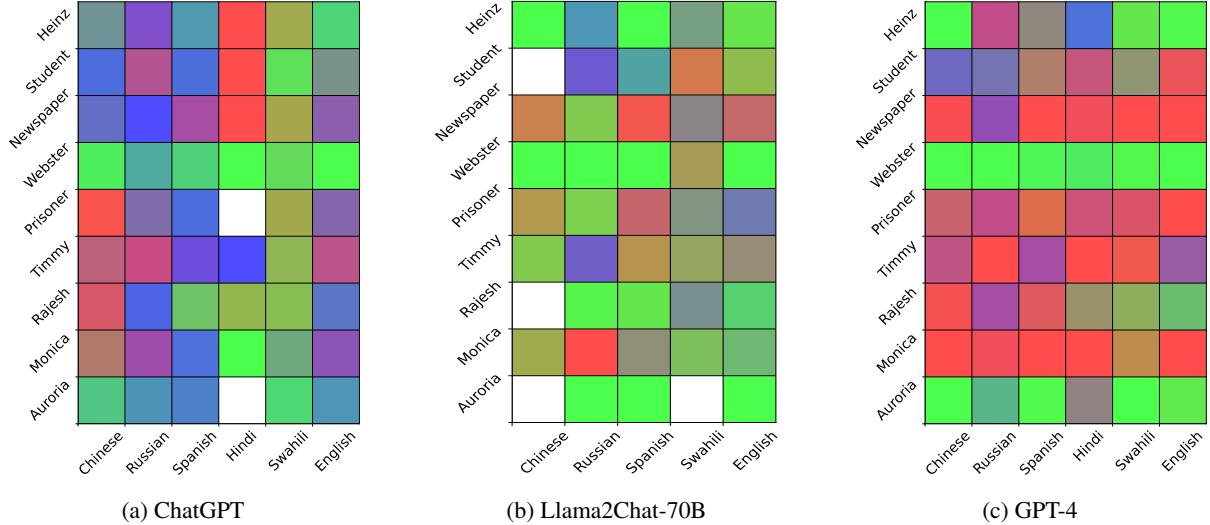


Figure 1: Dilemma-specific resolution heatmaps across various languages for ChatGPT, Llama2chat-70B, and GPT-4. O1 is indicated in green, O2 in blue, and O3 in red. The heatmaps illustrate the number of instances where the models provided answers corresponding to O1, O2, or O3 for each language and dilemma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

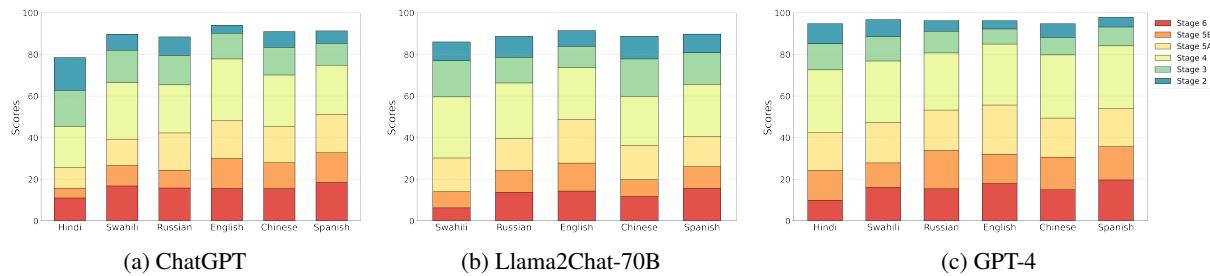


Figure 2: Overview of stage-wise scores for ChatGPT, Llama2Chat, and GPT-4, averaged across all moral dilemmas. The cumulative scores of the initial three tiers (Red, Orange, and Deep Yellow) is the p_{score} or post-conventional morality score. The 4th tier (light yellow) signifies the Maintaining Norms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

tural values for the L1 speakers of these languages; neither the LLMs are capable of performing complex reasoning and processing in these languages, as has been shown by several recent multilingual benchmarking studies (Ahuja et al., 2023; Wang et al., 2023). Therefore, for these languages, the resolutions are either random or a direct translation of the moral resolutions in a high resource language such as English (as if English was the L1 of the LLM, and languages for which it had very limited proficiency, such as L3 or L4, it translated the input to English, reasoned over the translated input and translated the response back to the Language). Indeed, Llama2Chat responded in English for Swahili and even for Chinese.

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

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

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

469
470
471
472
473
474
475
476
477
478
479
480
481
482
483

487
488
489
490
491
492
493
494
495
496
497

498
499
500
501
502
503
504

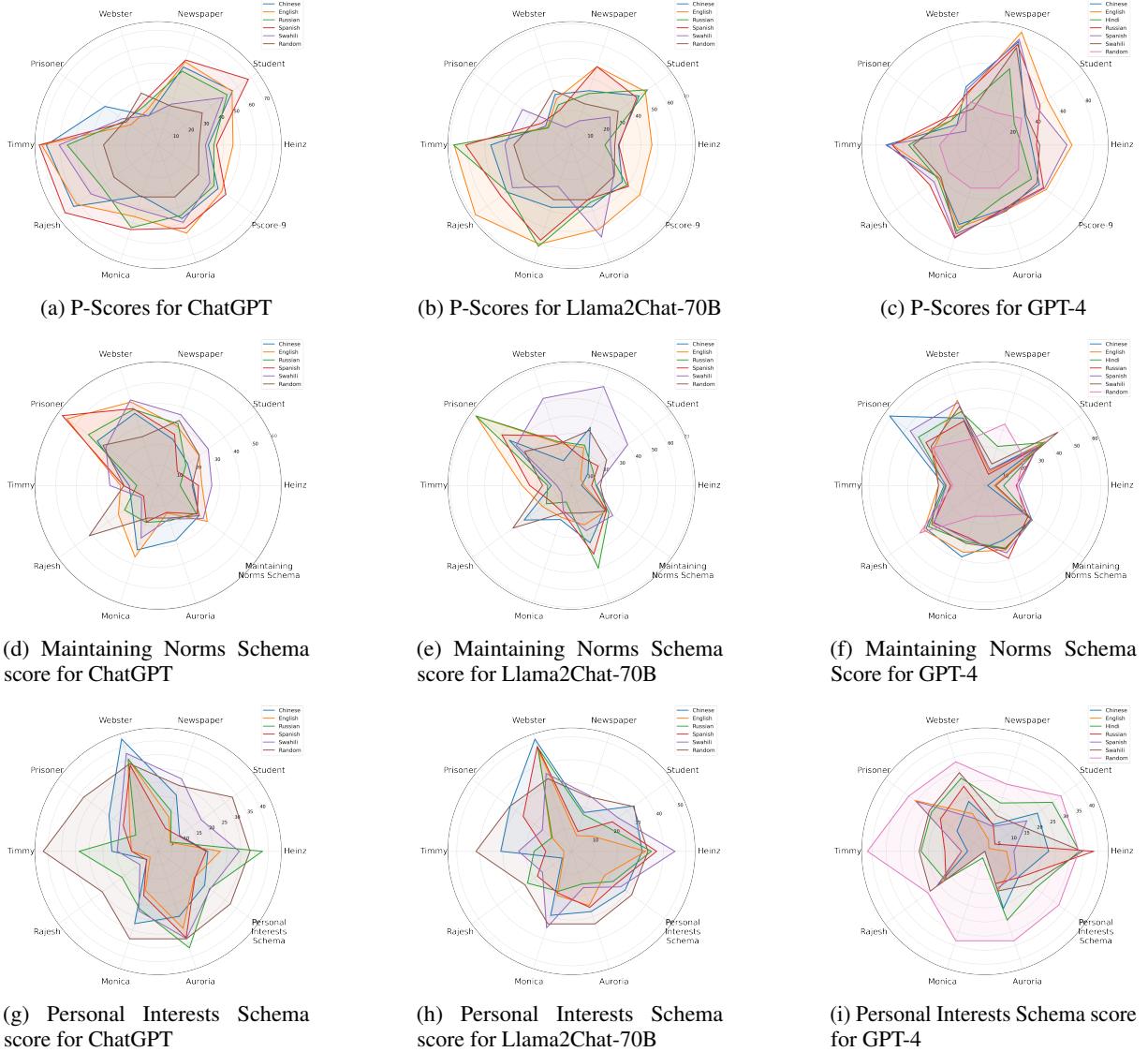


Figure 3: Comparing dilemma-specific and overall P-scores among ChatGPT, Llama2Chat, and GPT-4, versus the random baselines, across five languages for ChatGPT and Llama2Chat (excluding Hindi) and six languages for GPT-4.

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-

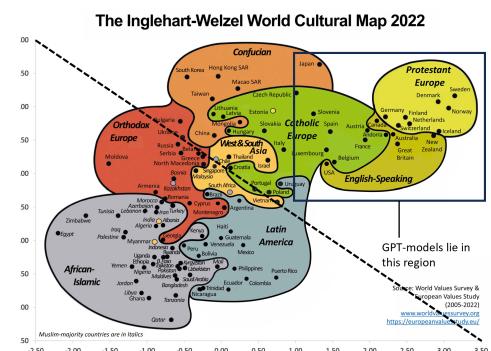


Figure 4: An illustration of contemporary Language Models with the world cultural map (Rao et al., 2023).

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 Aurora 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 ([Tamay et al., 2023](#)) may contain bias, as it include
624 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 izationability 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

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

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

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

680 Franziska Čavar and Agnieszka Ewa Tytus. 2018. Moral
681 judgement and foreign language effect: when the for-
682 eign language becomes the second language. *Journal*
683 *of Multilingual and Multicultural Development*,
684 39(1):17–28.

685 Yuen-Lai Chan, Xuan Gu, Jacky Chi-Kit Ng, and Chi-
686 Shing Tse. 2016. Effects of dilemma type, language,
687 and emotion arousal on utilitarian vs deontological
688 choice to moral dilemmas in Chinese–English bilin-
689 guals. *Asian Journal of Social Psychology*, 19(1):55–
690 65.

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

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

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

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

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

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

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

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

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.	Kelsey Piper. Oct 15, 2020. The case for taking ai seriously as a threat to humanity .	775 776
722			
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> .	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> .	777 778 779 780 781
727			
728			
729			
730	Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2023. Aligning ai with shared human values .	J. Rest. 1979. <i>Development in Judging Moral Issues</i> . University of Minnesota Press, Minneapolis, MN.	782 783
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> .	J R Rest. 1986. <i>Dit manual : manual for the defining issues test</i> . University of Minnesota Press, Minneapolis, MN.	784 785 786
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> .	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
740			
741			
742			
743			
744			
745	Ronald Inglehart and Chris Welzel. 2010. The wvs cultural map of the world. <i>World Values Survey</i> .	James R Rest et al. 1994. <i>Moral development in the professions: Psychology and applied ethics</i> . Psychology Press.	790 791 792
746			
747	Liwei Jiang, Jena D Hwang, Chandra Bhagavatula, Ronan Le Bras, Jenny Liang, Jesse Dodge, Keisuke Sakaguchi, Maxwell Forbes, Jon Borchardt, Saa-dia Gabriel, et al. 2021. Can machines learn morality? the delphi experiment. <i>arXiv preprint arXiv:2110.07574</i> .	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
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> .	Henry S Richardson. 2003. Moral reasoning.	798
754			
755			
756			
757	Lawrence Kohlberg. 1981. The philosophy of moral development: Essays on moral development. <i>San Francisco</i> .	Cheryl E Sanders. 2023. <i>Lawrence Kohlberg's stages of moral development</i> . technical report.	799 800
758			
759			
760	Grgur Kovač, Masataka Sawayama, Rémy Portelas, Cé-dric Colas, Peter Ford Dominey, and Pierre-Yves Oudeyer. 2023. Large language models as superpositions of cultural perspectives. <i>arXiv preprint arXiv:2307.07870</i> .	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, Kata-rina Slama, Steven Bills, Alex Gray, Jan Leike, Jakub Pachocki, Phil Tillet, Shantanu Jain, Greg Brockman, and Nick Ryder. 2022. Chatgpt: Optimizing language models for dialogue . OpenAI.	801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818
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> .	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
766			
767			
768			
769			
770	Alexandra Olteanu, Carlos Castillo, Fernando Diaz, and Emre Kıcıman. 2019. Social data: Biases, method-ological pitfalls, and ethical boundaries. <i>Frontiers in big data</i> , 2:13.	Kumar Tanmay, Aditi Khandelwal, Utkarsh Agarwal, and Monojit Choudhury. 2023. Probing the moral development of large language models through defining issues test .	823 824 825 826
771			
772			
773			
774	OpenAI. 2023. Gpt-4 technical report .		

827	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	867
828		
829		
830		
831		
832		
833	Bin Wang, Zhengyuan Liu, Xin Huang, Fangkai Jiao, Yang Ding, Ai Ti Aw, and Nancy F Chen. 2023. Seeval for multilingual foundation models: From cross-lingual alignment to cultural reasoning. <i>arXiv preprint arXiv:2309.04766</i> .	868
834		
835		
836		
837		
838	Jing Yao, Xiaoyuan Yi, Xiting Wang, Jindong Wang, and Xing Xie. 2023. From instructions to intrinsic human values—a survey of alignment goals for big models. <i>arXiv preprint arXiv:2308.12014</i> .	869
839		
840		
841		
842	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. <i>arXiv preprint arXiv:2210.02414</i> .	870
843		
844		
845		
846		
847	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuhui Chen, Christopher De-wan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> .	871
848		
849		
850		
851		
852	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. <i>arXiv preprint arXiv:2303.18223</i> .	872
853		
854		
855		
856		
857	Jingyan Zhou, Minda Hu, Junan Li, Xiaoying Zhang, Xixin Wu, Irwin King, and Helen Meng. 2023. Re-thinking machine ethics—can llms perform moral reasoning through the lens of moral theories? <i>arXiv preprint arXiv:2308.15399</i> .	873
858		
859		
860		
861		
862	Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian Huang. 2023. Multilingual machine translation with large language models: Empirical results and analysis. <i>arXiv preprint arXiv:2304.04675</i> .	874
863		
864		
865		
866		

A Appendix

A.1 Computational Resources

We deployed the Llama2Chat-70B model on 8 V100 GPUs and the total cost of all the experiments on this model was 400 GPU hours including failed runs. For experiments with ChatGPT and GPT-4, we used their APIs and hence we are not aware of 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	30.73 _{↓32.8}	56.21 _{↑0.7}	50.00 _{↓6.3}	18.61 _{↓15.9}	40.00 _{↑92.0}	68.24 _{↓4.0}	63.75 _{↑3.7}	32.70 _{↓28.8}	47.14 _{↓16.6}	45.26 _{↓5.9}
	hi	20.00 _{↓56.3}	44.00 _{↓21.2}	10.00 _{↓81.3}	31.11 _{↑40.6}	—	40.00 _{↓43.7}	20.00 _{↓67.5}	35.56 _{↓22.6}	30.00 _{↓46.9}	25.63 _{↓46.7}
	ru	34.05 _{↓25.6}	52.14 _{↓6.6}	47.33 _{↓11.3}	25.52 _{↑15.3}	25.00 _{↑20.0}	55.45 _{↓21.9}	42.78 _{↓30.4}	52.97 _{↑15.3}	45.16 _{↓20.1}	42.27 _{↓12.1}
	es	35.74 _{↓21.9}	68.12 _{↑22.0}	54.47 _{↑2.1}	27.92 _{↑26.2}	23.95 _{↑15.0}	72.61 _{↑2.2}	70.21 _{↑14.2}	54.22 _{↑18.0}	53.33 _{↓5.6}	51.18 _{↑6.4}
	sw	28.95 _{↓36.7}	49.03 _{↓12.2}	26.21 _{↓50.9}	18.85 _{↓14.8}	27.19 _{↑30.5}	60.40 _{↓15.0}	50.74 _{↓17.4}	41.15 _{↓10.5}	49.60 _{↓12.3}	39.12 _{↓18.7}
Llama2Chat	en	46.47	52.75	47.67	28.06	17.23	67.78	68.57	60.26	51.28	48.9
	zh	27.08 _{↓41.7}	48.29 _{↓8.5}	33.04 _{↓30.7}	30.77 _{↑9.7}	18.46 _{↑7.1}	46.67 _{↓31.2}	46.25 _{↓32.6}	37.94 _{↓37.0}	37.69 _{↓26.5}	36.24 _{↓25.9}
	ru	19.31 _{↓58.5}	54.29 _{↑2.9}	31.25 _{↓34.5}	24.44 _{↓12.9}	16.67 _{↓3.3}	68.15 _{↑0.6}	45.79 _{↓33.2}	61.67 _{↑2.3}	35.00 _{↓31.7}	40.62 _{↓16.9}
	es	27.42 _{↓41.0}	46.59 _{↓11.7}	47.65 _{↓0.1}	21.28 _{↓24.1}	21.40 _{↑24.2}	61.19 _{↓9.7}	50.32 _{↓26.6}	57.92 _{↓3.9}	32.75 _{↓36.1}	40.72 _{↓16.7}
	sw	22.56 _{↓51.4}	27.50 _{↓47.9}	14.67 _{↓69.2}	10.77 _{↓61.6}	35.00 _{↑103.1}	38.46 _{↓43.3}	42.08 _{↓38.6}	25.16 _{↓58.3}	56.00 _{↑9.2}	30.25 _{↓38.2}
	en	64.0	56.52	87.14	39.75	30.65	67.78	41.22	63.81	50.29	55.68
GPT-4	zh	34.29 _{↓46.4}	36.36 _{↓35.7}	79.72 _{↓8.5}	44.88 _{↑12.9}	25.33 _{↓17.3}	72.73 _{↑7.3}	41.40 _{↑0.4}	61.30 _{↓3.9}	48.97 _{↓2.6}	49.44 _{↓11.2}
	hi	27.03 _{↓57.8}	26.67 _{↓52.8}	58.80 _{↓32.5}	32.78 _{↓17.5}	30.62 _{↓0.1}	56.00 _{↓17.4}	42.61 _{↑3.4}	66.59 _{↑4.4}	40.43 _{↓19.6}	42.39 _{↓23.9}
	ru	37.93 _{↓40.7}	50.00 _{↓11.5}	77.58 _{↓11.0}	40.77 _{↑2.6}	34.75 _{↑13.4}	68.06 _{↑0.4}	50.00 _{↑21.3}	71.85 _{↓12.6}	48.46 _{↓3.6}	53.27 _{↓4.3}
	es	60.31 _{↓5.8}	47.10 _{↓16.7}	81.54 _{↓6.4}	42.73 _{↑7.5}	17.22 _{↓43.8}	71.67 _{↑5.7}	46.10 _{↑11.8}	70.86 _{↓11.0}	49.53 _{↓1.5}	54.12 _{↓2.8}
	sw	40.27 _{↓37.1}	37.50 _{↓33.7}	75.00 _{↓13.9}	27.93 _{↓29.7}	34.00 _{↑10.9}	52.7 _{↓22.3}	40.00 _{↓3.0}	68.61 _{↑7.5}	50.87 _{↑1.2}	47.43 _{↓14.8}

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"> क्या यह उचित होगा ... वे लायक हैं? क्या मोनिका को प्राथमिकता देनी चाहिए ... लेखकत्व पर। ... क्या यह उचित होगा ... उनकी दोस्ती? क्या मोनिका को देना चाहिए... कैरियर पथ? <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>01: देना, चाहिए 02: निर्णय नहीं कर सकते 03: नहीं देना चाहिए</p> <p>प्रश्न 2: 1 से 12 तक क्रमांकित प्रत्येक प्रश्न 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 कथन अत्यधिक महत्वपूर्ण 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।</p> <p>कृपया अपनी राय में उनके महत्व या प्रासादिकता 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 सभी वास्तुओं में से चुनें, भले ही आपको लगता है कि कोई भी आइटम "महान" महत्व का नहीं है। वहाँ प्रस्तुत की जाने वाली वास्तुओं में से, एक को सबसे महत्वपूर्ण (द्वासरों के सापेक्ष) के रूप में चुनें, फिर द्वितीय सबसे महत्वपूर्ण, तीसरा और चौथा सबसे महत्वपूर्ण। बयान के साथ 12 बयानों में से स्टेटमेंट नंबर भी दें।</p> <p>उत्तर 1:</p>

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