Is It Bad to Work All the Time? Cross-Cultural Evaluation of Social Norm Biases in GPT-4

Anonymous ACL submission

Abstract

LLMs have been demonstrated to align with the values of Western or North American cultures. 002 Prior work predominantly showed this effect through leveraging surveys that directly ask originally people and now also LLMs - about their values. However, it is hard to believe that LLMs would consistently apply those values 007 in real-world scenarios. To address that, we take a bottom-up approach, asking LLMs to reason about cultural norms in narratives from different cultures. We find that GPT-4 tends to generate norms that, while not necessarily incorrect, are significantly less culture-specific. 013 In addition, while it avoids overtly generating 014 stereotypes, the stereotypical representations of certain cultures are merely hidden rather than suppressed in the model, and such stereo-017 types can be easily recovered. Addressing these challenges is a crucial step towards developing 019 LLMs that fairly serve their diverse user base.

1 Introduction

022

024

LLMs are trained on vast web text. In principle, this data is representative of the diverse population of web users, which should contribute to LLMs serving the diverse population of their users. In practice, the training data predominantly consists of English web text from Western web users (Hershcovich et al., 2022), therefore covering more knowledge about Western cultures. Moreover, learning about the world through the lens of a Western user may entail that knowledge about other cultures is more prone to stereotyping and biases (Said, 1978; Nisbett et al., 2001; Henrich et al., 2010; Bender et al., 2021; Li et al., 2022).

Secondly, while web texts are authored by numerous web users, LLMs are trained on them as a single stream of unattributed text. As a result, they don't represent any specific person but rather an authoritative "voice from nowhere" which is supposedly representative of the diversity of its user



Figure 1: Top-down vs. bottom-up approaches to evaluating cultural alignment of LLMs. The top-down method asks direct survey-style questions about values, while the bottom-up approach asks models to reason about social norms in cultural narratives.

population but in practice is more aligned with a "default" user demographic (Cao et al., 2023a; Liu et al., 2023; Arora et al., 2022; Hartmann et al., 2023).

041

043

044

045

047

048

051

052

054

057

060

Finally, additional design choices in the development of LLMs, such as curation of (often proprietary) training data and "guardrails" designed to prevent models from generating harmful or stereotypical language, further leaks the values and norms of the (typically Western) developers into the models.

Prior work (Cao et al., 2023b; Ramezani and Xu, 2023; Durmus et al., 2024, inter alia) captured LLMs' cultural alignment and biases by leveraging existing surveys such as the Hofstede Culture Survey (HCS; Hofstede, 1984), the World Values Survey (WVS; Haerpfer et al., 2020), and the Global Attitudes Survey (PEW), finding that LLMs exhibited a strong alignment with North American cultures, and to a lesser extent, with other Western English speaking countries. However, such surveys that are designed to ask people about their values implicitly assume that people are consistent between their reported values and their real-life behavior. LLMs, on the other hand, don't have a consistent "persona" and are optimized to generate human-like responses.

In this work, rather than asking LLMs questions such as "For an average Chinese, how important it is to do work that is interesting (1-5)?", we embed these cultural aspects into narratives, such as the one presented in Figure 1. Such narratives may capture the more nuanced ways in which cultural conditioning implicitly affects people's everyday decisions and judgments (Selbst et al., 2019). Consequently, evaluating the responses from LLMs for such narratives can help identify inherent biases in LLM-backed decision making.

077

079

087

096

100

101

103

105

106

108

We adopt a bottom-up experimental design and use existing, human-written narratives from different cultures – specifically, plots from English Wikipedia for movies produced in various countries: China, India, Iran, and the United States. We instruct both annotators from the respective countries as well as GPT-4 to reason about the social norms in the movies, in the form of rules of thumb (RoT; Forbes et al., 2020). The annotators then judge the RoTs for their accuracy, culturespecificity, and stereotypicality.

We find that GPT-4 tends to generate norms that, while not necessarily incorrect, are significantly more generic. While GPT-4 generated norms were considered less stereotypical – likely thanks to its "guardrails"– reversing the question and asking GPT-4 to predict the agreement of people from certain countries with a particular norm resurfaces stereotypes and reveals the superficiality of the guardrails. Our study thus sheds light on the default representational biases of a prominent commercial language model, demonstrating that these models fail to live up to idealized (and probably impossible) egalitarian representations of a global public, and instead recapitulate the usual ever-present East vs. West racial hierarchies (Said, 1978).¹

Content Warning: This work contains examples that potentially implicate stereotypes, associations, and other harms that could be offensive to individuals in certain regions.

2 Background

Liu et al. (2024) define cultures using a taxonomy that includes cultural concepts, knowledge, values, norms and morals, linguistic form, and artifacts. In this paper, we focus on evaluating LLMs' judgments pertaining to **social norms** (§2.1), how well they align with various cultures, and to what extent these models are reinforcing **stereotypes** (§2.2).

2.1 Values, Norms, and Morals

With the recent progress in language technologies and their widespread adoption, there is vast interest in equipping these technologies with human-like values and norms.² Prior efforts in NLP focused on building norm banks for training norm-aware models (Forbes et al., 2020; Ziems et al., 2023), but they predominantly focused on Western norms (Liu et al., 2024).

At the same time, there is growing interest recently in serving users from diverse cultures (Hershowich et al., 2022). Various papers showed that LLMs exhibit a strong alignment with the values of North American cultures, and to a lesser extent, with other WEIRD countries, raising concerns about fairness (See for example, Johnson et al., 2022; Ramezani and Xu, 2023; Havaldar et al., 2023; Arora et al., 2023; Cao et al., 2023b; Santurkar et al., 2023; Tao et al., 2024; Durmus et al., 2024; Wang et al., 2024a; Masoud et al., 2025). Several of these papers experimented with different types of prompts, including mentioning the country name ("cultural prompting") or translating the prompt to the local language. These experiments typically reveal that when prompted in English without mentioning a cultural context, models by default assume a Western or even US culture.

The vast majority of studies in this area leverage existing surveys such as the Hofstede Culture Survey (HCS; Hofstede, 1984), which is centered around power distance, uncertainty avoidance, individualism-collectivism, masculinity-femininity, and short vs. long-term orientation; the World Values Survey (WVS; Haerpfer et al., 2020), which involves questions pertaining to social values, attitudes and stereotypes, well-being, trust, and more; 120

121

122

123

124

125

126

127

128

129

130

131

132

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

109

110

111

¹We will make data and code available upon publication.

²Liu et al. (2024) define *norms and morals* as a "set of rules or principles that govern people's behavior and everyday reasoning", making the distinction from *values*, which are defined as "beliefs, desirable end states or behaviors ranked by relative importance that can guide evaluations of things". We largely ignore this distinction in this paper.

| Country | Min | Max | Mean | Median |
|---------------|-----|-----|-------|--------|
| Token Count | | | | |
| China | 165 | 296 | 222.2 | 198.5 |
| India | 172 | 299 | 226.1 | 221.5 |
| Iran | 56 | 362 | 156.9 | 132.5 |
| United States | 171 | 276 | 217.4 | 224.0 |
| Verb Count | | | | |
| China | 16 | 40 | 27.80 | 28.5 |
| India | 18 | 49 | 31.95 | 31.0 |
| Iran | 5 | 53 | 21.40 | 17.0 |
| United States | 21 | 47 | 30.35 | 30.5 |

Table 1: Descriptive statistics of token counts and verb counts across movie plots by country.

or the Global Attitudes Survey (PEW),³ which asks people about their views on current global affairs. These studies present the survey questions to LLMs, directly asking them about their values.

153

155

156

157

158

159

160

161

163

164

165

167

168

170

171

172

173

174

175

176

178

179

181

182

183

184

185

186

189

With the caveat of social desirability bias (Grimm, 2010) and other factors which may affect people's responses, we can expect people to be largely consistent between their reported values and their real-life behavior. LLMs, on the other hand, don't have a consistent "persona" and are optimized to generate human-like responses. Thus, rather than asking questions about values and norms directly, Wang et al. (2024b) and Rao et al. (2025) start with prescribed social norms and use LLMs to generate more natural narratives in which these values should be considered. We rather take a *bottom-up* approach, prompting GPT-4 to reason about social norms in *existing* narratives from different cultural contexts.

2.2 Stereotypes and Cultural Bias

LLMs learn societal biases from their web-based training data, pertaining to race, gender, religion, profession, and more (Nadeem et al., 2021; Jha et al., 2023). Modern LLMs such as Gemini (Team et al., 2023) and GPT-4 (Achiam et al., 2023) do a better job at avoiding generating harmful or offensive content, thanks to their instruction tuning and preference tuning steps and other proprietary "guardrails" implemented by their developers. However, these superficial avoidance strategies likely only mask rather than remove the biases in these models. For example, Reuter and Schulze (2023) reveal the superficiality of the "guardrails" by showing that merely including the word "Muslim" in the prompt increased ChatGPT's response refusal rate - likely due to the association of this group with the hate speech it encounters online.

| Top Keywords | Торіс |
|---|---|
| rescue, captain, aircraft, bomb, ship film, village, children, doctor, women marriage, father, daughter, wife, Rajesh love, school, marry, proposes, life cow, teacher, village, barn, son friend, crush, high, student, picture girlfriend, baby, sister, father, love government, president, future, mother, son police, law, prison, duty, media home, husband, family, mother, house | military rural life family romance rural life romance family politics law enforcement family |
| | |

Table 2: Top keywords for each topic extracted using BERTopic from the movie plot dataset, along with our interpretation of the topic theme.

In another line of work, researchers revealed that in some setups, LLMs still express subtle or mild stereotypes towards various population groups, such as describing Arab characters as "poor and struggling" (Naous et al., 2024a) and Black people as "tall and athletic" (Cheng et al., 2023). This is especially concerning given the rise of popularity in using LLMs to generate synthetic users and study participants (Boelaert et al., 2025). 190

191

192

193

194

196

197

198

199

200

201

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

223

224

225

In this work, we contribute to this line of work by showing that when asked to generate cultural norms, GPT-4 avoids generating stereotypes. However, when used to simulate the agreement of people from that country with the stereotype, it predicts they would agree with it.

3 Data

In contrast to prior work that asks models about their values directly in a *top-down* approach, we take a *bottom-up* approach, presenting models with narratives from different cultural contexts and prompting them to reason about the social norms that these situations invoke. To that end, we first scraped the plots of movies produced in various countries from English Wikipedia, to serve as culturally-grounded narratives (§3.1). We then prompt annotators from the respective countries (§3.2), as well as GPT-4 (§3.3), to list the social norms invoked in these narratives.

3.1 Culturally-Grounded Narratives

To explore how social norms are reflected through culturally grounded narratives that affect people's everyday decisions and judgments, we focus on movie plots – widely consumed narrative media that often depict rich social behaviors and implicit norms. Such movie plots allow for context-rich interpretation and cultural priming.

³https://www.pewresearch.org/



Figure 2: Overview of RoT collection and evaluation process. We first scrape movie plots from English Wikipedia. Each plot is shown to both human annotators and GPT-4. Each human annotator (after the cultural activation task) write 3–5 RoTs per movie, while GPT-4 is prompted under two settings: default prompting and cultural prompting. This results in a total of 19–25 RoTs per movie across both human and model sources.

Following prior work (Shen et al., 2024a; Huang and Yang, 2023; Qiu et al., 2025), we focus on four geographically diverse countries: United States, China, India, and Iran. China and India were selected due to their large populations and globally recognized cultural distinctiveness; Iran represents a smaller culture group with unique traditions and perspectives; and the US represents a Western country overrepresented in English web text. We scraped the movies from English Wikipedia,⁴ retaining 20 movie plots for each country, and ensuring a moderate length to facilitate smooth annotation.⁵ Table 1 summarizes the dataset statistics.

To evaluate the diversity of social and cultural themes in our dataset, we applied topic modeling using BERTopic (Grootendorst, 2022). Table 2 presents representative keywords for each of the 10 topics identified by the model, along with our interpretation of the topic theme. We observe a wide range of themes, including family dynamics, romantic relationships, rural life, political events, and law enforcement. This thematic variety ensures that the cultural norms derived from the plots reflect a rich and heterogeneous set of lived experiences.

3.2 Human-written RoTs

Social Norms Format. To investigate culturally grounded social norms, we follow Forbes et al. (2020) and describe social norms in the Rules of Thumb (RoTs) format. RoTs are short declarative

statements describing appropriate or expected behavior. RoTs typically conform to the form "It is [judgment] [action]", where *judgment* is an adjective (e.g., "immoral" in Fig 2) and *action* is a clause (e.g., "to carry on an affair while being married").

255

256

257

258

259

260

261

263

264

265

266

267

268

269

270

272

273

274

275

276

277

278

279

280

281

282

283

286

287

Annotators. Figure 2 illustrates the RoT collection process which we detail below. We recruited annotators from the respective countries through the CloudConnect platform by Cloud Research.⁶ as well as through word of mouth. Annotators were compensated \$20-25 USD per hour. We collected annotations from 88 annotators across four cultural groups. Among those who reported, annotators ranged in age from 18 to 62 years (M = 34.1, SD = 9.5). The gender distribution was balanced, and the majority held a bachelor's degree. Detailed demographic breakdowns are provided in Appendix A.1.

Cultural Priming. Following Bhatia et al. (2024), to ensure their cultural affinity, we recruited annotators that have lived in the respective country for at least 5 years in the past 15 years. With that said, by design, CloudConnect annotators reside in English speaking countries, making them bicultural. Thus, to activate the cultural identity associated with the study (other than their current country of residence), we applied cultural priming, a technique widely validated in cultural psychological research (Hong et al., 2000; Oyserman and Lee, 2008; Liu et al., 2015). Specifically, annotators go through a small *cultural activation task* before their annotation task, in which they are shown five images pertaining to their culture, such as cultural icons, country flags, historical sites, and festivals, and are tasked with answering questions about the

⁴We chose English Wikipedia as it offers the most comprehensive and consistent coverage of international films in a single language, facilitating downstream analysis without requiring multilingual NLP tools.

⁵We randomly sampled 20 movies from the movies that fall between the 40th and 60th percentiles in terms of length for each country.

⁶https://www.cloudresearch.com/

<Country>-culture-driven: RoTs should align with established norms and practices in <Country>. Judgment + Action: Each RoT is in a single sentence with a straightforward structure: it is [the judgment] of [an action]. Verb-centric: Anchor each RoT to a specific verb from the story. Specificity: Avoid overly generic statements.

Table 3: Instructions for the RoT writing task, adapted from Forbes et al. (2020).

images to make sure they perceived and reflected on the priming material (See Appendix A.2).

290

294

302

305

307

313

314

315

317

321

323

324

327

RoT Writing Task. After completing the cultural activation task, annotators were asked to read a movie plot from their culture and provide 3–5 RoTs that are invoked by the narrative and that they perceive would be accepted within their culture. Each movie was annotated by three annotators. To help them come up with RoTs, we highlighted all the verb phrases in the plot as potential action terms, and prefilled a dropdown box with the 625 judgment adjectives from Forbes et al. (2020). See Table 3 for the annotation instructions and Appendix A.2 for the interface. Overall, we collected 396-441 human-written RoTs per culture.

3.3 GPT-generated RoTs

We use GPT-40 (OpenAI, 2024) with few-shot learning to generate 5 RoTs for a given movie plot. We prompted the model twice for each movie in the following setups:

Default Prompting: We ask the model to generate RoTs without referencing any cultural background (see Appendix B.1 for the prompt). This setup allows us to learn about the model's "default" cultural values.

Cultural Prompting: Mirroring the human annotation setup, we added "As someone with a <Country> cultural background..." to the default prompt. This framing encourages the model to generate culturally-aligned responses, simulating the perspective of a person from the specified country.

We acknowledge that some content in the movie plot – such as mention of cultural traditions, concepts, or names – may leak information about the culture to the model in the default prompting setup, making this setup less than 100% cultureagnostic. However, this setup provides less direct information about the target culture than the cultural prompting setup.

4 **Results**

We address three core research questions: (1) Which cultures does GPT-4 know about? We compare the accuracy and culture-specificity of GPT-4generated norms to human-written ones (§4.1). (2) Which cultures is GPT-4 aligned with? We measure which cultures align best with its default, culturally unmarked judgments (§4.2). (3) Does GPT-4 reinforce stereotypes? We measure the stereotypicality of GPT-4 generated norms and judgments (§4.3). 328

329

330

331

332

333

334

335

336

339

340

341

342

343

344

345

346

347

348

349

350

353

354

356

357

358

359

360

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

4.1 Which Cultures does GPT-4 Know about?

Human Ratings of RoTs. To assess the correctness and cultural alignment of the GPT-generated RoTs compared to the human-written ones, we recruited different sets of annotators from the target countries in a similar process to the annotation task in §3.2. Annotators similarly went through a cultural priming step before they proceeded to the main task. Each RoT was judged by 5 annotators – in the context of the plot from which is emerged – with respect to the following criteria (see Appendix A.2 for the annotation task):

- 1. Accuracy: On a Likert-scale from 1 to 5, to what extent the RoT accurately represents a social norm.
- 2. **Cultural-specificity:** On a Likert-scale from 1 to 5, to what extent the RoT reflects norms unique to the target culture, vs. a more generic or nearly-universal accepted norm.
- 3. **Stereotypicality:** Whether the RoT reinforces stereotypes about the target culture.

We collected annotations from 56 annotators across countries. Among those who reported, annotators ranged in age from 18 to 66 years (M = 33.3, SD = 10.6). The gender distribution was relatively balanced, and the majority held a bachelor's degree. Detailed demographic breakdowns are provided in Appendix A.1.

GPT-4-generated RoTs are – by and large – as accurate as human-written RoTs. Figure 3(A) presents the mean accuracy across RoTs for each combination of country and condition (humanwritten, GPT-4 generated with default prompting, and GPT-4 with cultural prompting). Overall, GPT-4- generated RoTs are rated as fairly accurate (≥ 3.8) across countries. The accuracy is identical to that of the human-written RoTs for India and Iran, and slightly higher (+0.1) than humanwritten RoTs for the US, but the difference is not



Figure 3: Average (A) accuracy, (B) cultural-specificity, and (C) stereotypicality scores for each country and condition combination.

statistically significant ($\beta = -0.31$, p = .003). For China, however, a small gap of +0.2 points favoring human-written RoTs was found to be statistically significant ($\beta = 0.25$, p = .001).⁷ There were no statistically significant differences between the two prompting strategies, indicating no clear advantage from cultural prompting, as was previously shown (Cao et al., 2023a). It's possible that the movie plot already provides implicit cultural cues, making it unnecessary to explicitly include the country name in the prompt.

377

378

384

394

400

401

402

403

404

405

406

407

408

409

410

GPT-4-generated RoTs are less culture-specific and more generic than human-written RoTs. Figure 3(B) shows the average specificity ratings for each combination of country and condition. Overall, RoTs were ranked as moderately culturespecific (1.7 – 2.8), suggesting a good number of generic or supposedly-universal RoTs across conditions. With that said, across countries, humanwritten RoTs were rated as more culturally-specific than GPT-4-generated RoTs ($\beta = 0.17, p = .037$) – suggesting that GPT-4's accuracy could in part be attributed to its tendency to generate generic norms that people across cultures can agree with. Again, cultural prompting showed no significant advantage over default prompting.

4.2 Which Cultures is GPT-4 Aligned with?

To evaluate which cultural perspective GPT-4 aligns with most closely, we reverse the roles and ask it to rate the accuracy of RoTs as a person from country X. The idea is that if GPT-4 shares implicit assumptions with a given culture, its default ratings should be closely aligned with those generated under that culture's perspective.



Figure 4: Average JSD between GPT-4's default predictions and culture-specific ratings across four countries. The model aligns most closely with the United States norms and deviates most from Chinese norms.

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

We randomly sampled 20 RoTs from each country and prompted GPT-4 to rate the accuracy of each RoT on a Likert scale of 1–5. We prompt the model in different settings; in the **default prompting** setting, we ask the model to rate the RoT without referencing any cultural background, testing its default, culture-agnostic stance. In the **cultural prompting** setting, we instruct the model to rate the RoT "As someone with a [country X] cultural background...", where X is each one of the countries in our experiments. For each setting, we estimate the distribution over ratings by prompting the model to generate 30 independent responses per RoT, using a temperature of 0.8 to introduce variability. See Appendix B.2 for the prompts.

Following Durmus et al. (2024), we use Jensen-Shannon Divergence (JSD) to measure distance between distributions. Specifically, we are interested in the distance between the distributions obtained from the default prompt and the country-specific distributions. We average the JSD values across the 20 RoTs to obtain an overall measure of cultural alignment for each culture.

⁷OLS regression: $F(11, 1645) = 30.69, p < .001, R^2 = 0.17.$

GPT-4 is mostly aligned with the US. As shown in Figure 4, the average divergence is lowest for the United States (0.12), indicating that GPT-4's default ratings are most similar to those given when prompted from the US perspective. Divergence increases for Iran (0.20), India (0.24), and is highest for China (0.33), suggesting that GPT-4's implicit normative stance deviates most from Chinese cultural framing. These results are consistent with prior findings that LLMs tend to default to dominant or Western cultural perspectives (Durmus et al., 2024; Naous et al., 2024b; Saha et al., 2025). The poorer performance on Iranian social norms compared to the US corroborates prior findings (Shen et al., 2024b; Saffari et al., 2025).

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

4.3 Does GPT-4 Reinforce Stereotypes?

Human-written RoTs are more stereotypical than GPT-4's. Figure 3(C) shows the average stereotypicality ratings from annotators across countries and conditions. While the majority of RoTs were perceived to be non-stereotypical, human-written RoTs were significantly more stereotypical than GPT-4-generated ones across countries ($\beta = 0.11, p < .001$). This seeminglysurprising finding could be explained by the following. First, one reason that the human-written RoTs are not completely stereotype-avoidant may be that people are less careful to avoid stereotyping their own group because there is more tolerance for stereotypes coming from in-group members (Bourhis et al., 1977; Thai et al., 2019). Another reason could be that GPT-4 lacks culture specific knowledge, including knowledge of stereotypes (Zhou et al., 2025), relative to the crowdworkers. Indeed, the Pearson correlation between culturalspecificity and stereotypicality across all samples yields a moderate positive correlation (r = 0.462, p < .001). Finally, GPT-4's stereotype avoidance could be attributed to its preference tuning and other safety mechanisms that prevent it from generating harmful, biased, or offensive responses. Indeed, previous work showed that LLMs are less likely to be critical of social groups (Boelaert et al., 2025; Zhang et al., 2025).

478GPT-4's biases are hidden, not removed.GPT-4794 appears to be stereotype-avoidant on the surface.480However, given the technical difficulty of debias-481ing, there are reasons to believe that the training482process that led to stereotype-avoidance is merely483hiding the bias rather than removing it – akin to

other ML models (Gonen and Goldberg, 2019; Kong et al., 2024). Modern language models are preference-tuned in the aims of avoiding generating text that are harmful towards minority groups (consider, for example, Hartvigsen et al., 2022). However, preference-tuning on text-generation datasets may fail to mitigate other forms of stereotypical judgments that the model can make, through classifications or ordinal predictions (Hofmann et al., 2024). That is, by shifting to a constrained ordinal setting, we hypothesize we can surface stereotypical associations between cultural groups and social norms that are prevalent in pretraining data (Caliskan et al., 2017). 484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

Specifically, to assess whether stereotypical associations can be recovered via a classificatory prompt, we start by identifying RoTs that were deemed as highly stereotypical by the human annotators. Then, we select the top quartile of RoTs with the highest stereotypicality scores within each cultural group. Additionally, we manually include a small number of RoTs whose negation reflects culturally salient stereotypes (e.g., "It is bad to work all of the time", China). This results in a set of 338 highly stereotypical norms (see Appendix C for examples).

We follow the same approach as in Sec. 4.2 to obtain the distribution of accuracy scores for each country as well as an unmarked (default) country. Figure 5 shows example RoTs where GPT-4's ratings appear to reflect cultural stereotypes. The first 5 distributions reflect the distribution of responses from GPT-4 and the last one (*human annotators*) shows the distribution of votes from annotators of the target culture.

On the top left, GPT-4 predicted that most Iranians would strongly agree that "It is not encouraged to live alone as a young bachelor", whereas human judgments about the accuracy of this RoT varied considerably. The predictions for China and India are similar, but in contrast, GPT-4 predicted that people from the US would strongly disagree with this RoT.

On the top right, GPT-4 predicted that people from India would strongly agree with the RoT "It is commendable to marry within your own religion" – accurately reflecting the human annotators from India but also perpetuating a norm that they judged as highly stereotypical. As harmless as this norm may seem, LLMs learning from web text might fail to account for changes in cultural norms over time (Zhou et al., 2025). In particular, while in-



Figure 5: Comparison of GPT-4 responses across four culturally grounded RoTs. Each subplot shows probability distributions of ratings (1–5) across five prompting conditions. Human annotations were provided by individuals from the cultural group associated with each RoT.

terfaith marriage in India is still uncommon, there is some level of acceptance towards it (Verma and Sukhramani, 2018).

536

539

541

543

545

546

547

548

551

554

558

561

565

A similar phenomenon can be demonstrated with anti-stereotypes. On the bottom left, Chinese annotators' agreement with the norm "It is bad to work all of the time" was moderate. But when asked to rate the RoT as a Chinese person, GPT-4 assigned the highest probability to "disagree", suggesting that it assumes most Chinese people are ok with working all the time. This is in line with the stereotypical view of Chinese culture as endorsing overwork which is likely prevalent in the model's training data. This perception of Chinese culture fails to account for both individual differences as well as changes in cultural norms over time, such as the recent push back among the younger population against the "996" work schedule (9 a.m. to 9 p.m., 6 days a week; Zhu and Peng, 2023). As a control, GPT-4 predicts that a "default", culturallyunmarked person would highly agree that it's bad to work all the time, and its predictions for other countries vary but are not as overwhelmingly disagreeing as the Chinese predictions.

A similar behavior is observed for the RoT "It is unethical to tell a lie to get benefits for yourself" (bottom right in Fig. 5); GPT-4's predicted distribution for Chinese raters places the highest likelihood on a neutral rating, suggesting uncertainty about whether such behavior is wrong – reflecting the stereotype that Chinese people are dishonest. This is an oversimplification of Chinese values that assess the morality of deception in light of its effects and the broader context in which it occurred, in contrast to the Western perception that dishonesty is always bad (Blum, 2007; Kwiatkowska, 2015). Again, as a control, GPT-4's ratings for other countries show stronger disapproval. 566

567

568

569

570

571

572

573

574

5 Conclusion

We show that GPT-4 exhibits default represen-575 tational biases when reasoning about culturally-576 grounded social norms. Specifically, its latent cul-577 tural representation aligns most closely with the US 578 and least with China, with India and Iran falling 579 in between. Moreover, while the model tends to 580 avoid generating overtly stereotypical language, 581 these stereotypes are still implicitly ingrained in the model and can be resurfaced - due to lack of 583 real technical solutions. Finally, our findings also 584 highlight a key tension in the design of culturally-585 competent LLMs, which on the one hand need to 586 possess culture-specific knowledge, while on the 587 other hand risk perpetuating stereotypes about the same cultures. Addressing these challenges is cru-589 cial given the diverse user base of LLMs and their widespread usage in downstream applications. 591

Limitations

592

593Scope.Our study uses countries as a proxy for594cultures, which is the most common proxy in NLP595research despite its limitations (Zhou et al., 2025).596Due to the cost of human annotations and API calls,597we focused on four geographically- and culturally-598diverse countries, and only evaluated GPT-4, which599we selected due to its popularity and wide reach.600Finally, due to the relatively small number of hu-601man annotators from each culture, we did not study602individual differences between annotators in this603study. Future work would need to cover a wider604range of cultures and models to draw a complete605picture of LLMs' default cultural representations.

Cultural Grounding. In this paper, we deviated 607 from the common practice to prompt LLMs directly about their values and instead prompted them to reason about social norms in existing narratives. We intentionally looked for human-written (as op-610 611 posed to LLM-generated) narratives grounded in different cultures. We chose movies because they 612 often reflect cultural norms (Rai et al., 2025). Yet, 613 it is possible that movies exhibit a certain "reporting bias" to depict more unusual events. Further-615 more, to factor out the effect of the multilingual 616 capabilities of GPT-4 on our study, we strictly lim-617 ited the experiments to English text.⁸ It is possible that a movie plot in English Wikipedia has been written from the perspective of a Western editor (Kumar, 2021). This setup, and the availability of crowdsourcing workers, also required us to employ 622 bicultural annotators - individuals who identified with the target culture but currently live in English-624 speaking countries - which could have impacted their judgments. We attempted to activate a specific cultural identify through cultural priming techniques. Nevertheless, even with our simplifying 628 assumptions, our study takes a step forward from quantifying LLMs' cultural alignment through surveys with direct question about values.

Ethical Considerations

633

634

637

638

Annotator Selection and Compensation. The study was conducted with the approval of our institute's Behavioral Research Ethics Board that reviewed the data collection procedures to ensure they posed no risk of harm to human participants. Annotators were compensated fairly according to

CloudResearch's compensation guidelines, which exceed local minimum wage standards. All annotation instructions explicitly directed participants to avoid including any personally identifiable information in their responses.

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

Screening for Harmful Content. Prior to human evaluation, we conducted a thorough review of the movie plots to screen for and remove any harmful or unsafe content. These steps were taken to ensure ethical compliance, participant safety, and data integrity throughout the study.

Using Country as a Cultural Proxy. We also acknowledge that cultural identity does not map neatly onto geographic or national boundaries, and that cultural variation exists at the individual level, shaped by personal history and experience. However, for the purposes of this study, we use country as a proxy for cultural grouping, consistent with prior work.

Inadvertent Stereotypes. We used culturally relevant images to prime annotators before norm generation and collected social norms rooted in specific cultural contexts. While our intention was to support cultural reflection, we acknowledge that both the images and the resulting norms may inadvertently reflect or reinforce cultural stereotypes.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Arnav Arora, Lucie-aimée Kaffee, and Isabelle Augenstein. 2023. Probing pre-trained language models for cross-cultural differences in values. In *Proceedings* of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP), pages 114–130, Dubrovnik, Croatia. Association for Computational Linguistics.
- Arnav Arora, Lucie-Aimée Kaffee, and Isabelle Augenstein. 2022. Probing pre-trained language models for cross-cultural differences in values.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT)*, pages 610–623.
- Mehar Bhatia, Sahithya Ravi, Aditya Chinchure, EunJeong Hwang, and Vered Shwartz. 2024. From

⁸In preliminary experiments we also tested translating prompts to the local language, which yielded subpar results.

- 698 702 712 713 714 715 717 719 720 721 723 724 725 726 727 728 729 730 731 732 734 735 737 738 739 740

- 741 742 743

local concepts to universals: Evaluating the multicultural understanding of vision-language models. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 6763–6782, Miami, Florida, USA. Association for Computational Linguistics.

- Susan D Blum. 2007. Lies that bind: Chinese truth, other truths. Rowman & Littlefield Publishers.
- Julien Boelaert, Samuel Coavoux, Étienne Ollion, Ivaylo Petev, and Patrick Präg. 2025. Machine bias. how do generative language models answer opinion polls?1. Sociological Methods & Research.
- Richard Y. Bourhis, Nicholas J. Gadfield, Howard Giles, and Henri Tajfel. 1977. Context and ethnic humour in intergroup relations. In ANTONY J. CHAPMAN and HUGH C. FOOT, editors, It's a Funny Thing, Humour, pages 261–265. Pergamon.
- Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. Science, 356(6334):183-186.
- Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. 2023a. Assessing cross-cultural alignment between chatgpt and human societies: An empirical study. arXiv preprint arXiv:2303.17466.
- Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. 2023b. Assessing cross-cultural alignment between ChatGPT and human societies: An empirical study. In Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP), pages 53-67, Dubrovnik, Croatia. Association for Computational Linguistics.
- Myra Cheng, Esin Durmus, and Dan Jurafsky. 2023. Marked personas: Using natural language prompts to measure stereotypes in language models. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1504-1532, Toronto, Canada. Association for Computational Linguistics.
- Esin Durmus, Karina Nguyen, Thomas Liao, Nicholas Schiefer, Amanda Askell, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, Liane Lovitt, Sam McCandlish, Orowa Sikder, Alex Tamkin, Janel Thamkul, Jared Kaplan, Jack Clark, and Deep Ganguli. 2024. Towards measuring the representation of subjective global opinions in language models. In First Conference on Language Modeling.
- Maxwell Forbes, Jena D. Hwang, Vered Shwartz, Maarten Sap, and Yejin Choi. 2020. Social chemistry 101: Learning to reason about social and moral norms. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 653-670, Online. Association for Computational Linguistics.

Hila Gonen and Yoav Goldberg. 2019. Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 609–614, Minneapolis, Minnesota. Association for Computational Linguistics.

744

745

747

748

749

752

753

754

755

756

757

758

760

761

762

763

764

765

766

767

768

769

777

778

779

780

781

782

783

784

785

788

789

790

791

792

794

795

796

797

798

799

- Pamela Grimm. 2010. Social desirability bias. Wiley international encyclopedia of marketing.
- Maarten Grootendorst. 2022. Bertopic: Neural topic modeling with a class-based tf-idf procedure. arXiv preprint arXiv:2203.05794.
- Christian Haerpfer, Ronald Inglehart, Alejandro Moreno, Christian Welzel, Kseniya Kizilova, Jaime Diez-Medrano, Marta Lagos, Pippa Norris, Eduard Ponarin, and Bi Puranen. 2020. World values survey wave 7 (2017-2020) cross-national data-set. (No Title).
- Jochen Hartmann, Jasper Schwenzow, and Maximilian Witte. 2023. The political ideology of conversational ai: Converging evidence on chatgpt's proenvironmental, left-libertarian orientation. arXiv preprint arXiv:2301.01768.
- Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. 2022. ToxiGen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3309–3326, Dublin, Ireland. Association for Computational Linguistics.
- Shreya Havaldar, Bhumika Singhal, Sunny Rai, Langchen Liu, Sharath Chandra Guntuku, and Lyle Ungar. 2023. Multilingual language models are not multicultural: A case study in emotion. In Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis, pages 202-214, Toronto, Canada. Association for Computational Linguistics.
- Joseph Henrich, Steven J Heine, and Ara Norenzayan. 2010. The weirdest people in the world? Behavioral and Brain Sciences, 33(2-3):61-83.
- Daniel Hershcovich, Stella Frank, Heather Lent, Miryam de Lhoneux, Mostafa Abdou, Stephanie Brandl, Emanuele Bugliarello, Laura Cabello Piqueras, Ilias Chalkidis, Ruixiang Cui, Constanza Fierro, Katerina Margatina, Phillip Rust, and Anders Søgaard. 2022. Challenges and strategies in crosscultural NLP. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 6997–7013, Dublin, Ireland. Association for Computational Linguistics.
- Valentin Hofmann, Pratyusha Ria Kalluri, Dan Jurafsky, and Sharese King. 2024. Ai generates covertly racist

910

911

- 801decisions about people based on their dialect. Nature,802633(8028):147–154.
 - Geert Hofstede. 1984. *Culture's consequences: International differences in work-related values*, volume 5. sage.
 - Ying-Yi Hong, Michael W Morris, Chi-yue Chiu, and Veronica Benet-Martínez. 2000. Multicultural minds:
 A dynamic constructivist approach to culture and cognition. *American Psychologist*, 55(7):709–720.
 - Jing Huang and Diyi Yang. 2023. Culturally aware natural language inference. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7591–7609. Association for Computational Linguistics.

810

811

813

814

815

816

817

819

823

824

825

830

831

834

837

839

843

845

847

849

- Akshita Jha, Aida Mostafazadeh Davani, Chandan K Reddy, Shachi Dave, Vinodkumar Prabhakaran, and Sunipa Dev. 2023. SeeGULL: A stereotype benchmark with broad geo-cultural coverage leveraging generative models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9851–9870, Toronto, Canada. Association for Computational Linguistics.
- Rebecca L Johnson, Giada Pistilli, Natalia Menédez-González, Leslye Denisse Dias Duran, Enrico Panai, Julija Kalpokiene, and Donald Jay Bertulfo. 2022. The ghost in the machine has an american accent: value conflict in gpt-3. *arXiv preprint arXiv:2203.07785*.
- Fanjie Kong, Shuai Yuan, Weituo Hao, and Ricardo Henao. 2024. Mitigating test-time bias for fair image retrieval. *Advances in Neural Information Processing Systems*, 36.
- Sangeet Kumar. 2021. *The digital frontier: Infrastructures of control on the global web.* Indiana University Press.
- Anna Kwiatkowska. 2015. How do others deceive? cultural aspects of lying and cheating. *The small and big deceptions: In psychology and evolutionary sciences perspective*, pages 46–72.
- Yizhi Li, Ge Zhang, Bohao Yang, Chenghua Lin, Shi Wang, Anton Ragni, and Jie Fu. 2022. Herb: Measuring hierarchical regional bias in pre-trained language models. In *Findings of the Association for Computational Linguistics: AACL-IJCNLP 2022*, pages 334–346. Association for Computational Linguistics.
- Chen Cecilia Liu, Iryna Gurevych, and Anna Korhonen. 2024. Culturally aware and adapted nlp: A taxonomy and a survey of the state of the art. *arXiv preprint arXiv:2406.03930*.
- Chen Cecilia Liu, Fajri Koto, Timothy Baldwin, and Iryna Gurevych. 2023. Are multilingual llms culturally-diverse reasoners? an investigation into multicultural proverbs and sayings.

- Zhuozhuo Liu, Menxue Cheng, Kaiping Peng, and Dan Zhang. 2015. Self-construal priming selectively modulates the scope of visual attention. *Frontiers in Psychology*, 6:1508.
- Reem Masoud, Ziquan Liu, Martin Ferianc, Philip C. Treleaven, and Miguel Rodrigues Rodrigues. 2025. Cultural alignment in large language models: An explanatory analysis based on hofstede's cultural dimensions. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 8474–8503, Abu Dhabi, UAE. Association for Computational Linguistics.
- Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. StereoSet: Measuring stereotypical bias in pretrained language models. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 5356–5371, Online. Association for Computational Linguistics.
- Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. 2024a. Having beer after prayer? measuring cultural bias in large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 16366–16393, Bangkok, Thailand. Association for Computational Linguistics.
- Tarek Naous, Michael J. Ryan, Alan Ritter, and Wei Xu. 2024b. Having beer after prayer? measuring cultural bias in large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1234–1245, Toronto, Canada. Association for Computational Linguistics.
- Richard E Nisbett, Kaiping Peng, Incheol Choi, and Ara Norenzayan. 2001. Culture and systems of thought: Holistic versus analytic cognition. *Psychological Review*, 108(2):291–310.
- OpenAI. 2024. Gpt-4o technical report. Accessed: 2025-05-15.
- Daphna Oyserman and Spike W. S. Lee. 2008. Does culture influence what and how we think? effects of priming individualism and collectivism. *Psychological Bulletin*, 134(2):311–342.
- Haoyi Qiu, Alexander Fabbri, Divyansh Agarwal, Kung-Hsiang Huang, Sarah Tan, Nanyun Peng, and Chien-Sheng Wu. 2025. Evaluating cultural and social awareness of llm web agents. In *Findings of the Association for Computational Linguistics: NAACL* 2025, pages 3978–4005, Albuquerque, New Mexico. Association for Computational Linguistics.
- Sunny Rai, Khushang Zaveri, Shreya Havaldar, Soumna Nema, Lyle Ungar, and Sharath Chandra Guntuku. 2025. Social norms in cinema: A cross-cultural analysis of shame, pride and prejudice. In *Proceedings* of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

967

968

913 914

912

915

916 917

919

921

928

931

934

935

936

937 938

941

942

943

944

947

949

951

953

955

957

958

959

960

961

962

963

964

965

966

Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 11396–11415, Albuquerque, New Mexico. Association for Computational Linguistics.

- Aida Ramezani and Yang Xu. 2023. Knowledge of cultural moral norms in large language models. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 428–446, Toronto, Canada. Association for Computational Linguistics.
- Abhinav Sukumar Rao, Akhila Yerukola, Vishwa Shah, Katharina Reinecke, and Maarten Sap. 2025. NormAd: A framework for measuring the cultural adaptability of large language models. In Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 2373–2403, Albuquerque, New Mexico. Association for Computational Linguistics.
 - Max Reuter and William Schulze. 2023. I'm afraid i can't do that: Predicting prompt refusal in blackbox generative language models. *arXiv preprint arXiv:2306.03423*.
 - Hamidreza Saffari, Mohammadamin Shafiei, Donya Rooein, Francesco Pierri, and Debora Nozza. 2025.
 Can I introduce my boyfriend to my grandmother? evaluating large language models capabilities on Iranian social norm classification. In *Findings of the Association for Computational Linguistics: NAACL* 2025, pages 6060–6074, Albuquerque, New Mexico. Association for Computational Linguistics.
 - Sougata Saha, Saurabh Kumar Pandey, and Monojit Choudhury. 2025. Meta-cultural competence: Climbing the right hill of cultural awareness. In Proceedings of the 2025 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 8025–8042, Albuquerque, New Mexico. Association for Computational Linguistics.
 - Edward W Said. 1978. Orientalism. Pantheon Books.
 - Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023.
 Whose opinions do language models reflect? In Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pages 29971–30004. PMLR.
 - Andrew D Selbst, Danah Boyd, Sorelle A Friedler, Suresh Venkatasubramanian, and Janet Vertesi. 2019.
 Fairness and abstraction in sociotechnical systems. In *Proceedings of the 2019 Conference on Fairness, Accountability, and Transparency*, pages 59–68. ACM.
 - Siqi Shen, Lajanugen Logeswaran, Moontae Lee, Honglak Lee, Soujanya Poria, and Rada Mihalcea.

2024a. Understanding the capabilities and limitations of large language models for cultural commonsense. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5668–5680, Mexico City, Mexico. Association for Computational Linguistics.

- Siqi Shen, Lajanugen Logeswaran, Moontae Lee, Honglak Lee, Soujanya Poria, and Rada Mihalcea. 2024b. Understanding the capabilities and limitations of large language models for cultural commonsense. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 5668–5680, Mexico City, Mexico. Association for Computational Linguistics.
- Yan Tao, Olga Viberg, Ryan S Baker, and René F Kizilcec. 2024. Cultural bias and cultural alignment of large language models. *PNAS Nexus*, 3(9):pgae346.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Michael Thai, Alex M. Borgella, and Melanie S. Sanchez. 2019. It's only funny if we say it: Disparagement humor is better received if it originates from a member of the group being disparaged. *Journal of Experimental Social Psychology*, 85:103838.
- Shweta Verma and Neelam Sukhramani. 2018. Interfaith marriages and negotiated spaces. *Society and Culture in South Asia*, 4(1):16–43.
- Wenxuan Wang, Wenxiang Jiao, Jingyuan Huang, Ruyi Dai, Jen-tse Huang, Zhaopeng Tu, and Michael Lyu. 2024a. Not all countries celebrate thanksgiving: On the cultural dominance in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6349–6384, Bangkok, Thailand. Association for Computational Linguistics.
- Yuhang Wang, Yanxu Zhu, Chao Kong, Shuyu Wei, Xiaoyuan Yi, Xing Xie, and Jitao Sang. 2024b. CDEval: A benchmark for measuring the cultural dimensions of large language models. In *Proceedings of the* 2nd Workshop on Cross-Cultural Considerations in NLP, pages 1–16, Bangkok, Thailand. Association for Computational Linguistics.
- Simone Zhang, Janet Xu, and AJ Alvero. 2025. Generative ai meets open-ended survey responses: Research participant use of ai and homogenization. *Sociological Methods & Research*, page 00491241251327130.
- Naitian Zhou, David Bamman, and Isaac L Bleaman. 2025. Culture is not trivia: Sociocultural theory for cultural nlp. *arXiv preprint arXiv:2502.12057*.

1024 1025 1026

- 1028

1029

- 1030 1031
- 1032

1036

1039

1040

1041

1043

1044

1045

1046

1047

1048

1049

1052

1053

1054 1055

1056

1057

1058

1059

1060

1061

1062

1064

Annotators Demographics A.1

А

Computational Linguistics.

Human Annotation

Images, 3(2):13–38.

RoTs Collection Task (§3.2). We collected annotations from 88 annotators across four cultural groups, primarily from the US (n = 44), followed by India (23), Iran (13), and China (10). Table 4 reports the country of residence composition for each culture. Among those who reported, annotators ranged in age from 18 to 62 years (M = 34.1, SD = 9.5), with most between 26 and 50. The gender distribution included 45 women (52.9%), 39 men (45.9%), and 1 non-binary participant (1.2%), and most reported holding a bachelor's (42.4%) or master's degree (24.7%).

Ying Zhu and Junqi Peng. 2023. From diaosi to sang

Caleb Ziems, Jane Dwivedi-Yu, Yi-Chia Wang, Alon

Halevy, and Divi Yang. 2023. NormBank: A knowl-

edge bank of situational social norms. In Proceed-

ings of the 61st Annual Meeting of the Association for

Computational Linguistics (Volume 1: Long Papers),

pages 7756–7776, Toronto, Canada. Association for

to tangping: The chinese dst youth subculture online. Global Storytelling: Journal of Digital and Moving

| Culture | Country of Residence | | | |
|---------|-----------------------------|----------------|---------------|--|
| | Canada | United Kingdom | United States | |
| China | 40.0 | 10.0 | 50.0 | |
| India | 13.0 | 8.7 | 78.3 | |
| Iran | 12.5 | 0.0 | 87.5 | |
| US | 0.0 | 0.0 | 100.0 | |

Table 4: Country of residence composition (%) of annotators within each cultural group for the RoT collection task.

RoTs Rating Task (§4). We collected annotations from 56 participants across four cultural groups, with the largest from the US (n = 25), followed by China (14), Iran (10), and India (7). Table 5 reports the country of residence composition for each culture. Participants ranged in age from 18 to 66 years (M = 33.3, SD = 10.6), with the majority between 18 and 50 years old. The sample included 28 women (50.0%), 19 men (33.9%), and 1 non-binary participant (1.8%), with most holding a bachelor's degree (50.0%).

A.2 Annotation Interface

Figures 6, 7, and 8 present the user interfaces for the cultural activation, RoT collection, and RoT evaluation tasks respectively, using Iran as an example culture.

| Culture | Country of Residence | | | |
|---------|----------------------|---------------|--|--|
| | Canada | United States | | |
| China | 28.6 | 71.4 | | |
| India | 16.7 | 83.3 | | |
| Iran | 0.0 | 100.0 | | |
| US | 0.0 | 100.0 | | |

Table 5: Country of residence composition (%) of annotators within each cultural group for the RoT evaluation task.

B LLM Prompts

RoT Collection B.1

Write rules-of-thumb to explain social norms to a computer. Your task is write social norms Rules-of-Thumb (RoTs) inspired by the provided story. I want RoTs that explain what we consider good and bad behaviour. You can think of the computer like a five-year-old child who doesn't know how the world works. We have to explain the basics to them about what people think is right and wrong in society, and how our actions make other people feel. Requirements for writing RoTs: Judgment+action: Each RoT is in a single 1. sentence with a straightforward structure: it is [the judgment] of [an action]. For example, It is rude to curse at people. 2. Verb-Centric: Anchor each RoT to a specific verb from the story (tense doesn't matter). For example: "curse; It is rude to curse at people." "likes; It is good to like your relatives." "invited; It is devastating to be excluded from a wedding you were invited to." 3. Specificity: Avoid overly generic statements. 4. Quantity: Write 3 to 5 RoTs. Bad Examples: 1. curse; Cursing at people is rude. (too rigid in format) 2. sleep; Studies have shown... (encyclopedic, not normative) 3. cut; It's ok to cut off contact... (requires original story context) 4. helped; It's kind to help people. (too vague)

Table 6: Prompt for the RoT collection task.

B.2 RoT Evaluation

1065

1066

1067

| Part 1. Answer the following questions. If unsure, provide the most likely answer. | | | | | |
|--|---|--|--|--|--|
| | Which country does this flag represent? Your answer | | | | |
| | What is the name of this figure? Your answer | | | | |
| | What is this place? Your answer | | | | |
| | What are the people doing? Use 1-2 sentences to describe about the thing they are making. Your answer | | | | |
| Happy Nourus! | What is this festival? Describe your experience with it in 2–3 sentences. Your answer | | | | |

Figure 6: Cultural activation task interface. Annotators are presented with five culturally relevant images (e.g., national flag, historical figures, landmarks, daily life, and festivals) and asked to answer short questions. This task primes participants to reflect on their cultural identity before writing social norms. Shown here is an example used for Iranian participants.

| su can think of the computer like a five-year-old child who doesn't know how the world works. We have to explain the basics to them about what people think is right and wrong in sciety, and how our actions make other people feel. | | | | | |
|--|---|--|---|--|--|
| Requirements for writing RoTs | | | | | |
| Clara Use - Cultu Sudgement + Act Verb-Centric: Anc © Verb-Centric: Anc © RoT1: Lcurs © RoT1: Lcurs © RoT1: Invite © RoT3: Invite © Specificity: Avoid © Quantity: Write 3 t | re-Driven: ion: Each F thor each F e -> It is ru g -> It's go e -> It's de overly gen to 5 RoTs. | RoTs should align with established norms and practices in Iran العراف . toT is in a single sentence with a straightforward structure: it is [the judgment] of [an action]. For example, It oT to a specific verb from the story (tense doesn't matter). For example: ide to curse at people. od to like your relatives. vastating to be excluded from a wedding you were invited to. aric statements. | is rude to curse at people. | | |
| 3ad Examples | | | | | |
| XRoT1: curse -> C Reason: It doesn't fol Correction: It is rude | ursing at p llow the str to curse at | eople is rude. ucture "It is [the judgment] of [an action]". people. | | | |
| XRoT2: sleep -> S Reason: RoTs are not Correction: It is ok to | tudies have t encyclope sleep for s | shown people perform best on exams after sleeping at least seven hours. dic knowledge. RoTs should contain everyday, commonsense knowledge about social norms and expectati even hours before exams. | ons. | | |
| KRoT3: cut -> It's Reason: An RoT mus Correction: It is cold | cold to cut t be fully u to cut off c | off contact with Jimmy. Iderstandable on its own, without the story it came from. Intact with your family member. | | | |
| KeoT4: helped -> Reason: Too vague. F Correction: It is kind t | lt's kind to RoTs should to help elde | help people. I be somewhat specific, balancing between explaining the underlying norms at play and applying to other si rly people who are struggling with standing on a bus. | tuations. | | |
| vu don't need to agree with t xw, please read the story first nvenience. | the above s st. Then, cl | ocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran المدان . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you ho although only having one hand, manages to sail his little boat. In his village, due to its hot climate and | can click the verb again to close it fo | | |
| u don't need to agree with t ow, please read the story first invenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. | the above s st. Then, cli is a sailor v ils are sen the country nts and st crewman. T | cocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran المدان . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you tho although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshic with his boat. At first he is reluctant, but because of the hardships of living he accepts the job. The crimin all the money needed for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he | I hard living conditions, J. Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first nvenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. | the above s st. Then, cl is a sailor v ils are <u>sen</u> the country nts and <u>st</u> crewman. T | ocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran المران . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you who although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshid with his boat. At first he is reluctant, but because of the hardships of living he accepts the job. The crimin all the money needed for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he | I hard living conditions, A. Khorshid is <u>asked</u> to illegally als <u>murder</u> one of the village's middle of the trip they <u>attack</u> himself <u>dies</u> due to the injuries | | |
| u don't need to agree with t w, please read the story first invenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. cape It is type a judgement | the above s st. Then, cl is a sailor v ils are <u>sen</u> the country nts and <u>st</u> crewman. T | cocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran الدان . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you tho although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshic with his boat. At first he is reluctant, but because of the hardships of living he accepts the job. The crimin all the money needed for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action | I hard living conditions, I hard living conditions, I. Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first nvenience. Captain Khorshidi dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. cape tt is type a judgement tt is type a judgement | the above s st. Then, cl is a sailor v ils are <u>sen</u> the countr- nts and <u>st</u> rrewman. T to to to | cocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran المدان . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the atthough only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshic with his boat. At first he is reluctant, but because of the hardships of living he accepts the job. The crimin all the money needed for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action type an action | I hard living conditions, I hard living conditions, I. Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first nvenience. Captain Khorshid i dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. It is type a judgement It is type a judgement It is type a judgement | the above s st. Then, cl is a sailor v is a sailor v sea sen the country nts and st crewman. T to to to to to | cocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran المران . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you tho although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshic with his boat. At first he is reluctant, but because of the hardships of living he accepts the job. The crimin al the money needed for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action type an action | I hard living conditions, I hard living conditions, I. Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first nvenience. Captain Khorshid i dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. It is type a judgement It is type a judgement It is type a judgement urder | the above s st. Then, cl is a sailor v lls are <u>sen</u> the countr the countr the countr the countr the countr the countr the countr to to to to | ocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran المران . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the although only having one hand, <u>manages</u> to <u>sail</u> his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they <u>ask</u> a middleman to <u>strike</u> a deal with Khorshic with his boat. At first he is reluctant, but because of the hardships of living he <u>accepts</u> the job. The crimin all the money <u>needed</u> for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is <u>killed</u> , Khorshid <u>faces</u> them single-handedly. He <u>manages</u> to <u>kill</u> all the criminals, but he type an action type an action | I hard living conditions, I hard living conditions, I. Khorshid is asked to illegally als <u>murder</u> one of the village's middle of the trip they <u>attack</u> himself <u>dies</u> due to the injuries | | |
| u don't need to agree with t w, please read the story first nvenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. It is type a judgement It is type a judgement It is type a judgement It is type a judgement It is type a judgement | the above s st. Then, cli is a sailor v sa sailor v second the countr. the countr. the countr. the countr. the countr. to to to to to | ocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran المران . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshid with his boat. At first he is reluctant, but because of the hardships of living he accepts the job. The crimin alt he money needed for the trip. At the beginning of the journey the criminals, kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action type an action type an action type an action | I hard living conditions, I hard living conditions, I Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first invenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. It is type a judgement It is type a judgement | the above s st. Then, cli is a sailor v ls are sen the countr, nts and st crewman. T to to to to to to to to to | cocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran المدان . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the atthough only having one hand, <u>manages</u> to <u>sail</u> his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshic with his boat. At first he is reluctant, but because of the hardships of living he accepts the job. The crimin all the money needed for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action type an action type an action type an action | I hard living conditions, I hard living conditions, I. Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first nvenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha khorshid and his of he sustained, it is type a judgement it is type a judgement | the above s st. Then, cl is a sailor v ils are sen the countr nts and st rrewman. T to to to to to to to to to to to to | cocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran المران . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshid rwith his boat. At first he is reluctant, but because of the hardships of living he accepts the niddleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action | I hard living conditions, J. Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first nvenience. Captain Khorshid i dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. t is type a judgement It is type a judgement | the above s st. Then, cl is a sailor v ils are <u>sen</u> the countr ths and <u>st</u> crewman. T to to to to to to to to | cocial norm RoTs. They're just examples. We want yours that reflect social norms in Iran المران . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshic with his boat. At first he is reluctant, but because of the hardships of living he accepts the piob. The criminal al the money needed for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action | I hard living conditions, I hard living conditions, I. Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first nvenience. Captain Khorshid i dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained, it is type a judgement It is type a judgement | the above s st. Then, cl is a sailor v lls are sen the countr the countr the countr the countr to to to to to to to to to to to | ocial norm RoTs. They're just examples. We want yours that reflect social norms in tran الدان . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the although only having one hand. manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshic with his boat. At first he is reluctant, but because of the hardships of living he accepts the job. The criminal the money needed for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action | I hard living conditions, I hard living conditions, I. Khorshid is asked to illegally als <u>murder</u> one of the village's middle of the trip they <u>attack</u> himself <u>dies</u> due to the injuries | | |
| u don't need to agree with t w, please read the story first nvenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha Khorshid and his of he isustained. It is type a judgement It is type a judgement | the above s st. Then, cli is a sailor v als are sen the countr the countr the countr the countr the countr to to to to to to to to to to to to to | ocial norm RoTs. They're just examples. We want yours that reflect social norms in tran الدان . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshid with his boat. At first he is reluctant, but because of the hardships of living he accepts the job. The criminal the money needed for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action | I hard living conditions, I hard living conditions, I. Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first invenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. It is type a judgement It is type a judgement | the above s st. Then, cl is a sailor v is a sailor v is a sent the countr, the countr, the countr, the countr, the countr, to | cocial norm RoTs. They're just examples. We want yours that reflect social norms in tran المران . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshid with his boat. At first he is reluctant, but because of the hardships of living he accepts the job. The crimin alt he money needed for the trip. At the beginning of the journey the criminals, kill the middleman, in the ne crewman is killed. Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action | I hard living conditions, I hard living conditions, I. Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first invenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained. It is type a judgement It is type a judgement | the above s st. Then, cli is a sailor v is a sailor v is a sailor v is a sent the countr, the countr, the countr, to | ocial norm RoTs. They're just examples. We want yours that reflect social norms in tran الجران Iran intro transmission in the area, so they ask a middleman to strike a deal with Khorshid into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshid with his boat. At first he is reluctant, but because of the hardships of living he accepts the lob. The criminal is the morey needed for the trip. At the beginning of the journey the criminals, kull the middleman, in the nerewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action type an action type an action | I hard living conditions, I hard living conditions, I. Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first invenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha Khorshid and his of he sustained; tis type a judgement it is type a judgement | the above s st. Then, cli is a sailor v sls are sen the countr, the countr, the countr, the countr, the countr, to | cocial norm RoTs. They're just examples. We want yours that reflect social norms in tran الجران . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshic with his boat. At first he is reluctant, but because of the hardships of living he accepts the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action | I hard living conditions, I hard living conditions, Khorshid is asked to illegally als murder one of the village's middle of the trip they lattack himself dies due to the injuries | | |
| u don't need to agree with t w, please read the story first invenience. Captain Khorshid dangerous crimina take them out of wealthiest mercha Khorshid and his c he sustained; tis type a judgement it is type a judgement | the above s st. Then, cl is a sailor v is a sailor v is are sen the countr the countr the countr to | cocial norm RoTs. They're just examples. We want yours that reflect social norms in tran الجران . ck on a verb and type judgements and actions. For each verb, you can write multiple RoTs. Once done, you the although only having one hand, manages to sail his little boat. In his village, due to its hot climate and into exile. They want to escape from the area, so they ask a middleman to strike a deal with Khorshic with his boat. At first he is reluctant, but because of the hardships of living he accepts the piob. The criminals all the money needed for the trip. At the beginning of the journey the criminals kill the middleman, in the he crewman is killed, Khorshid faces them single-handedly. He manages to kill all the criminals, but he type an action | I hard living conditions, I hard living conditions, Khorshid is asked to illegally als murder one of the village's middle of the trip they attack himself dies due to the injuries | | |

Figure 7: Interface for collecting Rules-of-Thumb (RoTs) from annotators. After completing the cultural activation task, annotators are shown a movie plot from their culture and asked to generate 3–5 RoTs that reflect social norms grounded in their cultural context. Action verbs in the story are automatically highlighted; clicking on a verb opens an input box where annotators can write norm statements anchored to that action. Instructions emphasize cultural alignment, verb specificity, and a clear judgment-action format (e.g., "It is rude to curse at people").

| Part 2. Need your help on evalu | uating "social norms" in Iran اليران! | |
|--|---|---|
| You will review several candidate "social norms" (highlight knowledge and experiences on the norms and practices i | ted in boxes). For each one, evaluate based on the following metrics. Keep in mind that your in Iran ايران . | responses should be based on your |
| Accuracy (rate 1-5): How accurately does it reflect Culture-Specificity (rate 1-5): How unique is this it Stereotype? (Yes or No): Does this statement rein than a nuanced or modern representation. | t the norms and practices in Iran الد ان based on your knowledge and experience? norm to Iran الدران A score of 5 means it is exclusive to this culture, while 1 indicates it is con force stereotypes about Iran الدران Consider whether it reflects an oversimplified or generalize | nmon across multiple cultures. ed perception of people in Iran ايران, rather |
| | | |
| | 93. It is admirable to face criminals alone | |
| | | |
| 1. To what extent is it accurate? Your Score: 3 | • | |
| Inaccurate | | cultur ایران Accurate to Iran |
| 2. To what extent is it unique to culture in Iran ایران ? Yo | bur Score: 3 | |
| Common across multiple cultures | • | cultur ایران Exclusive to Iran |
| 3. Does it reflect stereotypes about Iran ایران? | | |
| 0 | Yes, it's a stereotype. | |
| 0 | No, it's not a stereotype. | |
| | | |
| | 94. It is disgraceful to kill a child | |
| |) | |
| 1. To what extent is it accurate? Your Score: 3 | | |
| Inaccurate | • | Accurate to Iran ایران cultur |
| 2. To what extent is it unique to culture in Iran del? Yo | pur Score: 3 | |
| Common across multiple cultures | • | Exclusive to Iran العان cultur |
| | | |
| 3. Does it reflect stereotypes about Iran ایران? | | |
| 0 | Yes, it's a stereotype. | |
| 0 | No, it's not a stereotype. | |
| | | |
| | 95. It is tragic to kill someone in self-defense. | |
| | | |
| 1. To what extent is it accurate? Your Score: 3 | | |
| Inaccurate | | cultur ایران Accurate to Iran |
| 2. To what extent is it unique to culture in Iran ایران? Yo | bur Score: 3 | |
| Common across multiple cultures | • | ایران Exclusive to Iran ایران cultur |
| 3. Does it reflect stereotypes about Iran ایران? | | |
| 0 | Yes, it's a stereotype. | |
| 0 | No. it's not a stereotyne. | |
| | no, it a not a steleotype. | |

Γ

Figure 8: Evaluation interface. Annotators from each cultural background rate each RoT based on three criteria: accuracy, cultural specificity, and stereotypicality. They complete the same cultural priming task prior to rating, ensuring culturally grounded evaluations.

I will provide a statement (in double quotes). Estimate how people with a X cultural background would rate the extent to which this statement accurately reflects social norms and practices in X. Use a scale from 1 (not accurate) to 5 (accurate). Return only a single number from 1 to 5 and nothing else.

Table 7: Prompt for the RoT evaluation task. Red text was included only in the cultural prompting setup, with X set to one of the countries in our experiments (e.g., "Chinese background", "in China").

1068 C Example RoTs

1069Table 8 presents example stereotypical RoTs from1070each culture along with the human ratings for accu-1071racy, culture-specificity, and stereotypicality.

| Countr | y RoT | Accuracy | Specificity | Stereotypicality |
|--------|--|--|--|---|
| China | It is inappropriate to shut down a college. It is honorable to fulfill family obligations. It is important to study hard in school. It is important to study hard. It is good to be known for hard work. | [1, 2, 5] [5, 5, 5, 5, 3] [5, 5, 5, 5, 5] [5, 5, 5, 5, 5] [5, 5, 5] | $\begin{matrix} [4, 4, 3] \\ [4, 2, 4, 3, 1] \\ [3, 1, 4, 1, 5] \\ [3, 1, 3, 1, 5] \\ [3, 1, 3] \end{matrix}$ | $ \begin{bmatrix} 1, 1, 1 \\ 1, 1, 1, 1, 1 \\ 1, 1, 1, 1, 1 \end{bmatrix} $ |
| India | It is terrible to kill cows for human consumption. It is traditional to get an arranged marriage. It is mandatory to offer guests tea or coffee. It is dutiful to include all of your family members. It is responsible to arrange your sister's marriage. | [4, 5, 4, 2] [5, 5, 5, 5, 5] [5, 5, 5, 5, 2] [5, 5, 5, 5, 5] [5, 4, 5, 5, 5] | $\begin{bmatrix} 5, 5, 5, 3 \\ [4, 2, 2, 3, 5] \\ [2, 2, 3, 1, 5] \\ [5, 3, 2, 2, 5] \\ [4, 3, 2, 3, 5] \end{bmatrix}$ | $ \begin{bmatrix} 1, 1, 1, 1 \\ [1, 1, 1, 1, 1] \\ [1, 1, 1, 1, 1] \\ [1, 1, 1, 1, 1] \\ [1, 1, 1, 1, 0] \\ [1, 1, 1, 1, 0] \\ \end{bmatrix} $ |
| Iran | It is important for a woman to wear a chador outside. It is admirable to go the extra mile even when tired. It is rude to marry someone non-Iranian. It is immoral to reveal the body in public. It is okay to marry your cousin. | $ \begin{bmatrix} 5, 3, 4, 5, 1, 5 \end{bmatrix} \\ \begin{bmatrix} 4, 2, 4, 2, 4, 3 \end{bmatrix} \\ \begin{bmatrix} 3, 4, 4, 4, 4, 1 \end{bmatrix} \\ \begin{bmatrix} 5, 5, 4, 3, 3, 4 \end{bmatrix} \\ \begin{bmatrix} 4, 2, 4, 4, 5, 5 \end{bmatrix} $ | $ \begin{bmatrix} 5, 3, 4, 2, 5, 5 \end{bmatrix} \\ \begin{bmatrix} 4, 4, 4, 4, 4, 3 \end{bmatrix} \\ \begin{bmatrix} 3, 4, 4, 4, 4, 4 \end{bmatrix} \\ \begin{bmatrix} 5, 3, 4, 3, 3, 5 \end{bmatrix} \\ \begin{bmatrix} 4, 2, 4, 3, 5, 5 \end{bmatrix} $ | $ \begin{bmatrix} 0, 1, 1, 1, 1, 0 \\ [0, 0, 1, 1, 1, 1] \\ [1, 0, 1, 0, 1, 1] \\ [0, 1, 1, 0, 1, 1] \\ [0, 1, 1, 0, 1, 1] \\ [0, 1, 1, 0, 1, 1] \\ \end{bmatrix} $ |

Table 8: Selected stereotypical RoTs per culture, along with individual annotator ratings for accuracy, cultural specificity, and stereotypicality.