Global Gallery: The Fine Art of Painting Culture Portraits through Multilingual Instruction Tuning

Anonymous ACL submission

Abstract

Exploring the intersection of language and cul-001 ture in Large Language Models (LLMs), this study critically examines their capability to encapsulate cultural nuances across diverse linguistic landscapes. Central to our investigation are three research questions: the efficacy of 007 language-specific instruction tuning, the impact of pretraining on native language data, and the identification of optimal approaches to elicit accurate cultural knowledge from LLMs. Utilizing the GeoMLaMA benchmark for multilingual commonsense knowledge and an adapted 013 CAMeL dataset (english-only) for evaluation of nuanced cultural aspects, our experiments span six different languages and cultural con-015 texts, revealing the extent of LLMs' cultural 017 awareness. Our findings highlight a nuanced landscape: while language-specific tuning and bilingual pretraining enhance cultural under-019 standing in certain contexts, they also uncover inconsistencies and biases, particularly in non-Western cultures. This work not only expands our understanding of LLMs' cultural competence but also emphasizes the importance of integrating diverse cultural perspectives in their development, aiming for a more globally representative and equitable approach in language modeling.¹

1 Introduction

037

Large language models (LLMs) are capable of performing well across a wide variety of tasks (Bommasani et al., 2022; Srivastava et al., 2023) owing to their ability of generating coherent text that draws from a large corpus of pre-training data. However, some tasks like performing open-ended social reasoning involve questions (Parrish et al., 2022) which due to being under-specified or requiring a certain level of critical thinking elicit an opinionated answer from the LLM that affects different social groups, sometimes in undesirable ways (Bender et al., 2021). The role of culture is undeniable

¹Our code and data are available at this link



Figure 1: We instruction-tune LLaMA 2 in 5 non-English languages (Chinese, Hindi, Persian, Swahili, Greek) and evaluate both general cultural awareness as well as fine-grained cultural understanding in multilingual settings.

when looking at factors that determine people's beliefs and behavior in social settings. Cultural differences exist across countries and they interplay with the language spoken by native speakers, influencing both individual traits and group behavior. However, it is a well-known problem that multilingual LLMs are trained on corpora that are not equally representative of all parts of the world, but are rather more "western aligned" (Weidinger et al., 2022). This leads to potential issues of misrepresentation of culture and a lack of understanding of cultural knowledge in text generated by LLMs. Other work which studies this question brings out the lack of precision in cultural representations (Ramezani and Xu, 2023), problems of skewing distributions or amplifying biases existing in society (Jakesch et al., 2023), erasing underrepresented nuances (Hutchinson et al., 2020) and also the impact of low resourced languages and cultures they are

060

spoken in (Wibowo et al., 2023). While these studies have laid the groundwork, there remains a gap in understanding how language-specific instruction tuning might correlate with cultural knowledge.

061

062

063

067

068

094

100

101

102

103

104

106

107

108

109

110

Our work aims to first test the hypothesis that instruction tuning on data in a specific language might improve the cultural awareness of the LLM for the related culture. We also look at the impact of continued pre-training in a bilingual setting and how that has an influence on multilingual cultural understanding. Specifically we design the following research questions for this purpose:

> **RQ1 :** Does instruction tuning on languagespecific data enhance cultural knowledge?

> **RQ2a :** Does pretraining on language-specific data enhance cultural knowledge?

RQ2b: What is the optimal approach for eliciting cultural information from LLMs?

For this purpose, we translate the instructional data used for training the Alpaca model (Taori et al., 2023) to five languages (Chinese, Hindi, Persian, Swahili, Greek) other than English to cover the six distinct cultures. We then use this data to train low rank adapters (LoRA, Dettmers et al. (2023)) followed by evaluation on a benchmark for multilingual cultural knowledge in each language to measure the impact of instruction tuning.

We also explore whether LLMs understand tangible cultural nuances like food, beverages, clothing, etc by framing a dataset about different social situations with cultural targets based on a previous study (Naous et al., 2023). We use this to ask:

RQ3 : Do LLMs understand the nuances of culture and what disparities exist across tangible cultural aspects?

Overall, our findings show shortcomings in not only how culture is understood by LLMs, but also in current existing approaches at overcoming them.

2 Data and Methods

Our study explores the cultural understanding of large language models (LLMs) through two primary strategies: enhancement of an existing benchmark and the creation of a new, culturally-focused benchmark. Our methodology involves translating instructional data from the Alpaca dataset (which does not include culturally relevant information) into five additional languages and conducting supervised fine-tuning on various LLaMA 2 model (Touvron et al., 2023) sizes using these translations (a sample of which is manually verified by native speakers). With this approach we aim to test whether SFT in itself might improve cultural processing even though we are not explicitly training on culture specific data. This enables a broad examination of how LLMs handle cultural nuances across different linguistic contexts. 111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

2.1 Data

We work with 2 different datasets to understand cultural awareness at global and granular scales.

GeoMLaMA The GeoMLaMA benchmark (Yin et al., 2022) is central to our study on the cultural awareness of language models. Originally containing culturally diverse fill-in-the-blank sentences, we have converted it into a question-answer (QA) format. This adaptation makes it suitable for evaluating decoder-only models. Key features of this benchmark include:

- **Multilingual Scope:** Covers five countries (USA, China, India, Iran, Kenya), each with its native language (English, Chinese, Hindi, Persian, Swahili). We further expand our investigation by integrating a Greece/Greek variant of the dataset. This addition provides a broader spectrum for analysis, especially for languages that are lower resourced.
- **QA Format:** Consists of 900 multilingual questions with one gold correct and multiple incorrect answers, facilitating a clear assessment of the model's cultural understanding.
- **Cultural Diversity:** Questions cover a range of 17 broad cultural topics (eg. broom usage, climate, driver seat, measurement unit, etc) and are presented in both the country's native language and other languages, allowing for a comprehensive cross-cultural evaluation.

For **RQ1**, this dataset allows us to examine whether instruction tuning in a language specific to a given culture leads to better understanding and representation of that culture in language models. For **RQ2a** and **RQ2b**, the GeoMLaMA dataset's multilingual nature helps assess the impact of pretraining language models on language-specific data.

CAMeL Our study also incorporates the CAMeL dataset, initially introduced in the "Beer After Prayer" paper (Naous et al., 2023), to conduct a more detailed cultural analysis. Originally designed to compare Arabic and Western cultural norms, we have adapted the CAMeL dataset to align with the countries featured in the GeoM-LaMA benchmark. This adaptation involves col-

lecting new data from native speakers and select-161 ing sentence templates that are broadly applicable 162 across various cultures. Our modified version of 163 the CAMeL dataset is tailored to specifically ad-164 dress RQ3, which focuses on the language models' 165 granular understanding of cultural elements. 166

Key aspects of our adapted CAMeL dataset are: 167

168

169

170

172

173

174

175

176

177

178

179

180

181

182

183

185

186

187

188

189

191

- Cultural Adaptation: We have enriched CAMeL to reflect six cultures from the GeoMLaMA benchmark, involving data collection from native speakers and culturally diverse sentence templates.
 - Cultural Categories and Prompts: The dataset contains nine categories, such as gendered pairs, each with ten unique prompts and around fifty targets, covering a range of cultural elements like food, names, clothing, and literature.
 - QA Scenarios for Granular Analysis: We create five types of multiple-choice QA scenarios from CAMeL, designed to assess the models' depth of cultural understanding and their ability to distinguish between various cultural elements.

In our experiments, we deploy various multiplechoice QA scenarios derived from both datasets, ensuring that choice order is randomized to mitigate positional bias in large language models (Pezeshkpour and Hruschka, 2023). This approach allows us to comprehensively address each research question, ensuring that our findings are robust and well-supported by empirical evidence.

Language-specific finetuning data 2.2

To investigate the effect of language-specific instruction tuning on cultural awareness, we begin 193 with the 52k instruction-following demonstrations 194 used for training the Alpaca model (Taori et al., 195 2023), referred to as the cleaned Alpaca dataset. 196 These instructions, originally in English, are trans-197 lated into six languages (English, Chinese, Hindi, 198 Persian, Swahili, and Greek) using an automatic 199 translation system from the NLLB project (Team et al., 2022). These translations correspond to the 201 native languages of six cultures (American, Chi-202 nese, Indian, Iranian, Kenyan, and Greek) under study, resulting in the Alpaca-X dataset, where 'X' denotes the respective language for eg., Alpaca-en is the original English Alpaca data, and Alpaca-hi 206 is the translated Hindi version. It is important to 207 note that all these datasets are content-equivalent, only differing in terms of language.

2.3 Supervised Instruction Finetuning

We employ 4-bit QLoRA (Dettmers et al., 2023) to 211 train using supervised finetuning (SFT), low-rank 212 adapters (LoRA) for the base models using our 213 Alpaca-X data, with hyperparameters detailed in 214 the Appendix Table 12. These adapters, specific 215 to each language, can be integrated into the base 216 model in a plug-and-play manner. The base model 217 combined with a language-specific adapter trained 218 on Alpaca-X data is also referred to as an Alpaca-X 219 model, for simplicity of notation. For eg., Alpacahi data is used to train an adapter for the Alpaca-hi 221 model. Note that SFT, supervised finetuning, and 222 instruction tuning is used interchangeably through-223 out the rest of the document. 224

210

220

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

Experimental Settings 3

We divide our experiments into two distinct categories - first looking at how instruction tuning and pretraining play a role in cultural understanding, and then going deeper into different aspects of cultural nuances.

Studying the effects of language specific 3.1 instruction tuning

Our experiments are designed to isolate the impact of different components (base model, LoRA, evaluation prompt) on the cultural awareness of LLMs.

Experimental Setups To address our research questions, we have devised the following experimental setups:

- 1. For RQ1 (Language-Specific Instruction Tuning): We compare the LLaMA 2 model with an English-specific adapter (Alpaca-en) against Alpaca-X models, where 'X' denotes other languages. This comparison helps determine the effectiveness of language-specific instruction tuning in enhancing cultural understanding.
- 2. For RQ2a (Language-Specific Pretraining): We explore the performance of bilingual base models for Chinese $(Yi)^2$ and Swahili (Uliza)³, each with its respective LoRA, to gauge the impact of language-specific pretraining on cultural knowledge.
- 3. For RO2b (Quality of Fine-Tuning Data): An ablation study contrasts a non-Alpaca-X adapter, developed from high-quality bilingual data, with our Swahili Alpaca adapter. This helps assess

²Huggingface link for Yi

³Huggingface link for Uliza (pre-trained model)

the influence of fine-tuning data quality on cul-tural understanding.

258

259

261

262

263

266

270

271

272

274

277

281

283

287

290

298

Distribution of token counts for pretraining and instruction tuning In the context of our experiments, the token counts for pre-training and instruction tuning vary significantly. The Alpaca-X models are the only components developed in-house, while the pretrained bilingual models and Swahili LoRA are sourced from open-source repositories.

The LLaMA 2 model underwent pre-training with a substantial 2 trillion tokens. For the Swahili base model, a continued pre-training phase incorporated 0.32 billion Swahili tokens. In contrast, the Chinese base model involved pre-training with a combined total of 3 trillion Chinese and English tokens. The Alpaca dataset used for instruction tuning is relatively small, consisting of 52,000 instructions, which translates to approximately 5.817 million tokens for the English Alpaca, or 0.005 billion tokens. The token count for the non-Alpaca LoRA, used in one of our ablation studies, remains unknown because the open source repository it is adapted from does not specify details. This disparity in token counts highlights the differences in data scale between pre-training and instruction tuning phases which might have some effect on our results that we cannot control.

> Further, note that all these controlled experiments are performed using the first dataset (GeoMLaMA) without going into granular details, because our focus was to study the effects of instruction tuning instead of different cultural aspects.

3.2 Granular Analysis of Cultural Aspects

In this section, we outline a series of experiments utilizing the CAMeL dataset to conduct an in-depth analysis of cultural aspects. These experiments are designed to evaluate the model's nuanced understanding of cultural elements.

Setting 1: A multiple-choice question is framed with one option representing the answer from the corresponding culture and five options from the other cultures (all other settings are restricted to four options). This setting aims to assess the model's comprehension of individual cultural aspects rather than a general overview.

Setting 2: No options from the correct culture are
provided in the multiple-choice questions, options
are randomly sampled uniformly from incorrect
cultures. This approach is intended to determine
the model's default cultural inclination when the

Base Model	LoRA	Prompt
English	{lang} Alpaca	{lang}
English	English Alpaca	English
{lang}	{lang} Alpaca	{lang}
{lang}	Non-Alpaca LoRA	{lang}

Table 1: The four experimental combinations we test for RQ1 and RQ2. *lang* refers to language-specific variants of Alpaca or a language-specific prompt, translated from English.

correct option is absent.

Setting 3: Three options from the correct culture are provided alongside one option from a randomly selected incorrect culture. A model with accurate cultural understanding should consistently avoid choosing the incorrect option. This setting tests the model's alignment with the findings from the previous settings. 306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

Setting 4: Each question includes four options from the correct culture, but three of these are from a different category than what the question addresses to test for precision of understanding. For instance, in a question about names, three options might be food items, with only one being a name. The model's ability to discern between categories within the same culture is evaluated here. If the model understands culture minutely enough to be able to differentiate between the categories we are asking about, then it would never pick an option from the incorrect category. But if it only has a fuzzy understanding of culture, then it might end up choosing any of the given options as all of them are "culturally correct" in a global sense.

Setting 5: Questions regarding gendered categories are used, where half of the options are correct for the gender but incorrect for the culture, and the other half are correct for the culture but incorrect for the gender. This setting tests whether the model prioritizes cultural accuracy over gender accuracy in its responses. For instance, in a question about American female names with options as Liam, David, Aisha and Divya, we expect the model to choose one between Aisha and Divya over the two male names. Ideally, the model response should stick to the correct gender, because typically female clothing is worn by females and male clothing by males, and similarly usually females have female-associated names and vice versa. But if the model responses stick to the correct culture and ignore gender, then the model does not necessarily understand the details of the gendered cultural aspect even if it is broadly culturally correct.

348

351

352

364

366

370

371

374

375

378

379

381

394

3.3 Evaluation technique

Our evaluation method draws inspiration from existing approaches that aggregate token logprobabilities for prompt completion. Specifically, the techniques employed by Trinh and Le (2019) and Wang et al. (2023b) utilize variations of this concept. These studies use aggregated token logprobabilities in determining the most likely prompt completions.

Further building on this approach, we utilize CAPPr (K. Dubey, 2023), a tool that implements the aforementioned idea by selecting the completion most likely to follow a given prompt. CAPPr achieves this by calculating the log-probability of each token in a completion, considering both the prompt tokens and the preceding tokens within that completion. This process involves averaging the log-probabilities to derive the inverse perplexity of the completion. Subsequently, these averaged log-probabilities are exponentiated to obtain a completion probability. This procedure is repeated for each potential completion to form a normalized probability distribution over the set of completions, which for our use case represents the different options in a QA setting.

4 Results

Our findings show that LLMs do *not* understand the specific details that define culture even when we try different approaches like SFT, bilingual pretraining, and prompting in the native language.

4.1 RQ1 : Does SFT on language specific data enhance cultural knowledge? (No)

Our investigation into whether SFT on languagespecific data enhances cultural knowledge involves a series of experiments, detailed in Appendix A, Table 6. This section focuses on key results pertinent to our hypotheses for the RQ.

Better than BERT, but only in English Analysis of the GeoMLaMA performance (Table 2) compares the English Base Model combined with Language-Specific Alpaca and Language-Specific prompts against the English Base Model with English Alpaca and English prompts across three different model sizes ranging from 7B to 70B. We restrict our analysis to non-USA countries where English is not the native language. Results indicate that instruction tuning in English slightly outperforms encoder models like BERT and XLMR, a trend not always observed when SFT is applied



Figure 2: The 70B LLaMA 2 model shows strong performances for China and Iran across cultural concepts for different cultures.

SFT lang	China	India	Iran	Kenya	Greece
R	esults fron	n GeoML	aMA bei	nchmark	
(mBERT)	0.30	0.41	0.21	0.30	-
(XLMR-L)	0.37	0.37	0.37	0.32	-
	Prom	ot langua	ge: engl	ish	
eng (7)	0.50	0.39	0.24	0.31	0.34
eng (13)	0.54	0.42	0.31	0.28	0.34
eng (70)	0.46	0.45	0.28	0.28	0.38
	Prom	pt langua	ge: {lan	g}	
{lang} (7)	0.25	0.39	0.31	0.31	0.28
$\{lang\}(13)$	0.32	0.36	0.28	0.34	0.28
{lang} (70)	0.39	0.33	0.14	0.34	0.34

Table 2: Instruction tuning on language specific data does not consistently enhance cultural knowledge across languages and cultures. The numbers 7, 13 and 70 correspond to the model sizes in billions of parameters. The metric is the GeoMLaMA benchmark metric on a scale of 0-1 with higher being better.

in other languages. This discrepancy may be attributed to the predominance of English in the LLaMA 2's pretraining data, making it the language that is most coherent for the model. We also note that while some amount of cultural data is definitely present in the pretraining data, our SFT instructional data does not include cultural content. 397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

Larger model does not necessarily mean better cultural knowledge Contrary to expectations, larger models (7B, 13B, 70B parameters) do not consistently demonstrate enhanced cultural knowledge. This suggests that factors such as pretraining data composition or pre-training hyperparameters or even translation of the instruction tuning data from English, significantly influence the model's cultural understanding.



Figure 3: The distribution of countries chosen by the 70B LLaMA 2 model without the question explicitly mentioning the chosen country shows a large percentage favouring China and Iran.

No clear enhancement due to SFT, with a lot 413 of variability across cultures The hypothesis 414 415 that SFT with language-specific data substantially improves cultural knowledge is not conclusively 416 supported by our findings. There is notable vari-417 ability across different cultures. For instance, in 418 China and Iran, English-based fine-tuning seems to 419 420 be more effective, while in India, Hindi fine-tuning competes closely. English emerges as the most ef-421 fective language for eliciting cultural knowledge 422 across various cultures. However, the second-most 423 effective language is not consistently the native one. 494 For example, prompting in Chinese yields better 425 426 results than Swahili for Kenyan cultural questions (Table 6). This observation might be influenced by 427 the larger representation of Chinese compared to 428 Swahili in the pretraining data set. 429

4.2 RQ2a : Does pretraining on native language data enhance cultural knowledge? (Yes)

430

431

432

433

434

435

436

437

438

Our study also probes the influence of pretraining language distribution on cultural understanding. Specifically, we contrast the LLaMA 2 model, primarily pretrained on English data, with bilingual base models for Chinese (Yi) and Swahili (Uliza). The performance comparisons are in Table 3.

Pretraining is useful in improving cultural un-439 derstanding along with the instruction tuning 440 and language specific prompting In the con-441 text of English language queries, both Yi and 449 Uliza models do not surpass the performance of 443 LLaMA 2. However, for queries related to China, 444 when prompted in Chinese, the Yi model demon-445 strates superior performance compared to LLaMA 446 2 and also achieves parity with LLaMA 2's En-447

Model	Size	China	Kenya
Prompt langu	age : Er	ıglish	
	7	0.50	0.31
LLaMA 2 + eng Alpaca	13	0.54	0.28
	70	0.46	0.28
Vi Lang Almana	6	0.43	-
11 + eng Alpaca	34	0.39	-
Uliza + eng Alpaca	7	-	0.25
Uliza + {swa, eng} LoRA	7	-	0.31
Prompt language .	Chine:	se/Swahili	
	7	0.25	0.31
LLaMA 2 + zh/swa Alpaca	13	0.32	0.34
	70	0.39	0.34
Vi h Alassa	6	0.39	-
11 + Zn Alpaca	34	0.54	-
Uliza + swa Alpaca	7	-	0.31
Uliza + {swa, eng} LoRA	7	-	0.41

Table 3: Pretraining on language specific data helps to improve cultural awareness. Bilingual non-alpaca finetuning along with bilingual continually pretrained model gives the most culturally appropriate responses when prompted in the respective native language.

glish performance. Similarly, for Kenyan cultural queries, the Uliza model, when prompted in Swahili, matches LLaMA 2's performance. 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

The quality of pretraining data matters for increased awareness across cultures Notably, the Yi model generally outperforms LLaMA 2 in English for cultures outside China (USA, Iran, Kenya, Greece), as shown in Table 6. This suggests that high-quality, filtered pretraining data, particularly when used for continued pre-training, play a crucial role in enhancing a model's cultural awareness across different cultures.

4.3 RQ2b: Optimal Approach for Eliciting Cultural Knowledge

In an ablation study focusing on the quality of finetuning data, we examine a non-Alpaca LoRA derived from carefully curated Swahili data.⁴ This is contrasted with the performance of our Swahili Alpaca model. The results indicate a clear superiority of the curated LoRA over Alpaca, as it not only surpasses Alpaca but also outperforms the English results from LLaMA 2.

This finding underscores that the most effective approach for eliciting accurate cultural knowledge involves a bilingual base model pre-trained on high-quality, language-specific data. Additionally, supplementing this model with a LoRA,

⁴Huggingface link for Uliza (finetuned model)

instruction-tuned on curated instructional examples
and prompted in the respective language, further
enhances its cultural understanding. Such a combination of high-quality pretraining, targeted instruction tuning, and language-specific prompting
emerges as the optimal strategy for achieving deep
cultural insight.

4.4 RQ3 : Do LLMs understand granular tangible cultural aspects? (somewhat)

482

483

484

485

486

488

489

490

491

492

493

494

495

496

497

498

499

501

502

503

504

505

506

507

509

510

512

513

514

515

516

517

518

519

520

521

522

This part of our study, centered on English, aims to delve into the nuanced cultural understanding of models, using the CAMeL dataset. We chose English for several reasons: the complexity of translating proper nouns, the redundancy of translated nouns representing the same concept, and previous findings indicating superior English performance unless a language-specific pretrained model is available.

Prior distributions of cultural aspects of countries affect cultural understanding at the granular level Our analysis reveals that the model displays a pronounced preference for certain cultures, particularly China and Iran, when no correct options are present (Setting 2, Figure 3). This inherent bias significantly affects performance for other cultures when correct options are included.

As illustrated in Figure 2, LLaMA 2 generally exhibits the highest performance across various cultural aspects for China (Setting 1), with some exceptions where Iran leads. However, despite previous research indicating alignment with American values, the model shows a relatively superficial understanding of American culture, as evidenced by its lower performance. A possible explanation could be that the other options provided for the question have higher prior distributions, but there are possibly multiple factors at play here.

In another test (Setting 3), we present a scenario where three options are from the correct culture, along with one option from a randomly selected, incorrect culture. This setup is intended to evaluate the model's ability to discern cultural appropriateness accurately. Our findings reveal a stark contrast in performance based on the cultural context. For questions pertaining to China, the model demonstrates a high degree of accuracy, rarely selecting the incorrect cultural option. In contrast, when presented with questions about male names in the US, the model's performance significantly declines, choosing the incorrect option nearly 70%

Category	USA	China	India	Iran	Kenya	Greece
Beverage	0.56	0.31	0.67	0.40	0.50	0.63
Female Clothing	0.60	0.69	0.58	0.81	0.79	0.69
Female Names	0.89	0.87	0.97	0.82	0.92	0.85
Food	0.32	0.40	0.76	0.32	0.69	0.28
Literature	0.21	0.33	0.45	0.20	0.34	0.65
Location	0.81	0.88	0.76	0.72	0.84	0.81
Male Clothing	0.58	0.54	0.85	0.86	0.75	0.74
Male Names	0.94	0.85	0.97	0.85	0.93	0.87
Religion	0.51	0.55	0.81	0.72	0.53	0.66
Overall	0.60	0.61	0.76	0.65	0.70	0.69

Table 4: We measure the percentage of times that LLaMA 2 70B prefers an option from an incorrect category when provided with a single choice from the correct category paired with 3 incorrect ones. Ideally, this should be close to 0 if the model has true understanding.

of the time. This disparity highlights the model's uneven capability in distinguishing between culturally relevant and irrelevant options across different cultural settings. An alternate plausible explanation could be that model has learned that the US is a largely multi-cultural society, which our evaluation approach isn't designed to consider.

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

Such biases could be attributed to the mixed cultural perspectives inherent in the pretraining data, which might emphasize certain cultures over others. Detailed results across all culture dimensions and model sizes are provided in the Appendix, Tables 7, 8, and 9.

Complex cultural understanding is lacking even for countries for which the model understands culture broadly Our investigation further explores the model's depth of cultural understanding through a specific testing approach. In this setup, each question offers four culturally appropriate options, but only one option is relevant to the question's category, while the other three belong to different categories. The assessment focuses on the frequency with which the model selects an option from an incorrect category. Notably, even for countries like China and Iran, where the model generally shows a good grasp of broader cultural aspects, the selection of incorrect category options is alarmingly high, as detailed in Table 4. Ideally, the model should have a near-zero selection rate of incorrect categories for countries with strong cultural representation. However, this nuanced understanding appears to be lacking.

A striking example involves the 13B model's in-

terpretation of Chinese female names. In an array 558 of approximately 3K questions, the model consis-559 tently showed a preference for beverage names over 560 actual female names. For instance, it judged 'Goji Berry tea' as a more probable name than 'Chen' in the given context of 'I met a girl named [fill 563 in] at the park'. While it's conceivable that the 564 model might not always err when presented with non-beverage incorrect options, the fundamental issue remains that it should not select such incon-567 gruent options at all. Comprehensive results of this 568 testing are available in the Appendix, Table 10, un-569 derscoring the model's limitations in distinguishing between specific cultural categories. 571

Being culturally accurate is preferred by LLaMA 2 over being gender accurate, even though it should be the opposite In a nuanced test, we presented options that juxtapose cultural accuracy against gender accuracy: two options cor-576 rect in culture but incorrect in gender, and two others correct in gender but incorrect in culture. The 578 results reveal a marked preference for cultural accuracy over gender accuracy, particularly in contexts where the cultural representation in the model's training data is more pronounced (see Table 5 and 582 Appendix 11).

573

574

579

581

584

585

586

590 591

596

597

604

This tendency is more evident in countries with a higher cultural prominence in the model's training data. For instance, in questions related to China, the model predominantly selects culturally accurate responses, regardless of gender correctness. Conversely, for countries like the USA, the model shows a greater propensity to choose options that are correct in terms of gender. This pattern suggests that the prominence of certain cultural or gender concepts in the pretraining corpus along with grammatical gender signals in language significantly influences the model's decision-making process, underscoring the impact of training data composition on the model's understanding of nuanced cultural and gender-related aspects.

5 **Related Work**

In the context of understanding cultural biases in Large Language Models (LLMs), several studies have made significant contributions, each addressing different aspects of this multifaceted issue. Tao et al. (2023) employ the World Values Survey to map GPT models on the Inglehart-Welzel Cultural Map, highlighting the effectiveness of cultural prompting as a mitigation strategy. Durmus

Category	USA	China	India	Iran	Kenya	Greece
Female Clothing	0.40	0.83	0.59	0.31	0.44	0.54
Male Clothing	0.51	0.66	0.77	0.42	0.48	0.68
Female Names	0.14	0.90	0.58	0.61	0.66	0.58
Male Names	0.23	0.93	0.57	0.65	0.68	0.65
Overall	0.32	0.89	0.63	0.51	0.56	0.61

Table 5: LLaMA 2 70B prefers being culturally correct than being gender correct across cultures.

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

et al. (2023) combine datasets from the World Values Survey and the Pew Research Center's Global Attitudes surveys to explore models' alignment with Western values, using various prompting techniques. The SeaEval benchmark (Wang et al., 2023a) demonstrates the challenges multilingual LLMs face in multicultural reasoning, affected by factors like positional bias and language nuances. COPAL-ID (Wibowo et al., 2023) finds that LLMs have a lower understanding of culture-related questions compared to non-culture related ones, especially in multilingual settings. Additionally, Cao et al. (2023) pioneered examining cultural alignment from a chatbot perspective, revealing Chat-GPT's American-centric alignment.

However, these studies collectively highlight some gaps: a predominant focus on Westerncentric perspectives, limited exploration of non-Western cultures, and the need for more comprehensive strategies to incorporate a global spectrum of cultural nuances in LLMs. Our work aims to build upon these findings, addressing these shortcomings by examining LLMs' cultural awareness more holistically and inclusively.

6 Conclusion

This study on the cultural understanding of Large Language Models (LLMs) reveals significant variations in their ability to encapsulate diverse cultural nuances. Our investigations, leveraging the GeoMLaMA benchmark and the adapted CAMeL dataset, demonstrate that while language-specific instruction tuning and bilingual pretraining offer some improvements, they fall short of ensuring comprehensive cultural competence, particularly in non-Western contexts. The findings underscore the need for incorporating a wider range of cultural perspectives in LLM training and development, highlighting the importance of creating models that are not only linguistically adept but also culturally sensitive and globally inclusive.

7 Limitations

651

663

664

665

668

671

672

675

676

680

687

688

694

695

698

This study, while extensive, is subject to certain limitations which are important to acknowledge:

- 1. The current methodology conceptualizes culture as a singular entity within a nation-state. This perspective, while useful for structured analysis, might not fully capture the rich diversity and complexity of modern societies, where multiple cultures and languages coexist within a single country. Future research could benefit from exploring more granular approaches that can effectively address this multifaceted nature of cultural identity.
 - 2. The pretraining process lacks control over token distribution, contrasting with the controlled instructional data used in fine-tuning experiments. This could affect result interpretation. Future work should investigate the effects of smaller, high-quality datasets for controlled pre-training across languages.
 - Our experiments use 4-bit QLoRA for instruction tuning, and it's uncertain if results would differ with higher-bit configurations. Further research is needed to explore the impact of varying bit settings.
 - Evaluating large language models is an ongoing challenge within the field, and the methodology chosen for this study, while grounded in established research, has its strengths and limitations. This approach needs to be considered alongside alternative evaluation methods, each with their respective advantages and drawbacks, to suit specific use cases and research objectives.

References

- 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 '21, page 610–623, New York, NY, USA. Association for Computing Machinery.
- Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto,

Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. 2022. On the opportunities and risks of foundation models.

700

701

702

703

704

707

708

709

710

711

712

713

714

715

717

718

719

720

721

722

723

724

725

726

727

728

730

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

- Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. 2023. 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.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms.
- Esin Durmus, Karina Nyugen, Thomas I. 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. 2023. Towards measuring the representation of subjective global opinions in language models.
- Ben Hutchinson, Vinodkumar Prabhakaran, Emily Denton, Kellie Webster, Yu Zhong, and Stephen Denuyl. 2020. Social biases in NLP models as barriers for persons with disabilities. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 5491–5501, Online. Association for Computational Linguistics.
- Maurice Jakesch, Advait Bhat, Daniel Buschek, Lior Zalmanson, and Mor Naaman. 2023. Co-writing with opinionated language models affects users' views. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23, New York, NY, USA. Association for Computing Machinery.
- K. Dubey. 2023. Cappr (version 0.8.7) [computer software].

Tarek Naous, Michael J. Ryan, Alan Ritter, and Wei Xu. 2023. Having beer after prayer? measuring cultural bias in large language models.

759

765

774

777

779

793

796

810

811

812

813

814

815

818

- Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. 2022. BBQ: A hand-built bias benchmark for question answering. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2086–2105, Dublin, Ireland. Association for Computational Linguistics.
- Pouya Pezeshkpour and Estevam Hruschka. 2023. Large language models sensitivity to the order of options in multiple-choice questions.
- Aida Ramezani and Yang Xu. 2023. Knowledge of cultural moral norms in large language models.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Divi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia,

Fatemeh Siar, Fernando Martínez-Plumed, Francesca 819 Happé, Francois Chollet, Frieda Rong, Gaurav 820 Mishra, Genta Indra Winata, Gerard de Melo, Ger-821 mán Kruszewski, Giambattista Parascandolo, Gior-822 gio Mariani, Gloria Wang, Gonzalo Jaimovitch-823 López, Gregor Betz, Guy Gur-Ari, Hana Galijase-824 vic, Hannah Kim, Hannah Rashkin, Hannaneh Ha-825 jishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, 826 Hinrich Schütze, Hiromu Yakura, Hongming Zhang, 827 Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, 828 Jack Geissinger, Jackson Kernion, Jacob Hilton, Jae-829 hoon Lee, Jaime Fernández Fisac, James B. Simon, 830 James Koppel, James Zheng, James Zou, Jan Kocoń, 831 Jana Thompson, Janelle Wingfield, Jared Kaplan, 832 Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, 833 Jason Wei, Jason Yosinski, Jekaterina Novikova, 834 Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen 835 Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Ji-836 aming Song, Jillian Tang, Joan Waweru, John Bur-837 den, John Miller, John U. Balis, Jonathan Batchelder, 838 Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose 839 Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, 840 Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, 841 Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl 842 Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, 843 Katja Markert, Kaustubh D. Dhole, Kevin Gim-844 pel, Kevin Omondi, Kory Mathewson, Kristen Chi-845 afullo, Ksenia Shkaruta, Kumar Shridhar, Kyle Mc-846 Donell, Kyle Richardson, Laria Reynolds, Leo Gao, 847 Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-848 Ochando, Louis-Philippe Morency, Luca Moschella, 849 Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng 850 He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem 851 Şenel, Maarten Bosma, Maarten Sap, Maartje ter 852 Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas 853 Mazeika, Marco Baturan, Marco Marelli, Marco 854 Maru, Maria Jose Ramírez Quintana, Marie Tolkiehn, 855 Mario Giulianelli, Martha Lewis, Martin Potthast, 856 Matthew L. Leavitt, Matthias Hagen, Mátyás Schu-857 bert, Medina Orduna Baitemirova, Melody Arnaud, 858 Melvin McElrath, Michael A. Yee, Michael Co-859 hen, Michael Gu, Michael Ivanitskiy, Michael Star-860 ritt, Michael Strube, Michał Swedrowski, Michele 861 Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike 862 Cain, Mimee Xu, Mirac Suzgun, Mitch Walker, 863 Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor 864 Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun 865 Peng, Nathan A. Chi, Nayeon Lee, Neta Gur-Ari 866 Krakover, Nicholas Cameron, Nicholas Roberts, 867 Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas 868 Deckers, Niklas Muennighoff, Nitish Shirish Keskar, 869 Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan 870 Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, 871 Omer Levy, Owain Evans, Pablo Antonio Moreno 872 Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, 873 Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, 874 Percy Liang, Peter Chang, Peter Eckersley, Phu Mon 875 Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, 876 Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing 877 Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta 878 Rudolph, Raefer Gabriel, Rahel Habacker, Ramon 879 Risco, Raphaël Millière, Rhythm Garg, Richard 880 Barnes, Rif A. Saurous, Riku Arakawa, Robbe 881

Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima, Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Piantadosi, Stuart M. Shieber, Summer Misherghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models.

900

901

902

903

904

905

907

909

910

911

912

913

914

915

916

917

918

919

921

922

924

926

931

935

936

937

938

941

- Yan Tao, Olga Viberg, Ryan S. Baker, and Rene F. Kizilcec. 2023. Auditing and mitigating cultural bias in llms.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https:// github.com/tatsu-lab/stanford_alpaca.
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang.

2022. No language left behind: Scaling humancentered machine translation.

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models.
- Trieu H. Trinh and Quoc V. Le. 2019. A simple method for commonsense reasoning.
- Bin Wang, Zhengyuan Liu, Xin Huang, Fangkai Jiao, Yang Ding, Ai Ti Aw, and Nancy F Chen. 2023a. Seaeval for multilingual foundation models: From cross-lingual alignment to cultural reasoning. *arXiv preprint arXiv:2309.04766*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023b. Self-consistency improves chain of thought reasoning in language models.
- Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, Courtney Biles, Sasha Brown, Zac Kenton, Will Hawkins, Tom Stepleton, Abeba Birhane, Lisa Anne Hendricks, Laura Rimell, William Isaac, Julia Haas, Sean Legassick, Geoffrey Irving, and Iason Gabriel. 2022. Taxonomy of risks posed by language models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '22, page 214–229, New York, NY, USA. Association for Computing Machinery.
- Haryo Akbarianto Wibowo, Erland Hilman Fuadi, Made Nindyatama Nityasya, Radityo Eko Prasojo, and Alham Fikri Aji. 2023. Copal-id: Indonesian language reasoning with local culture and nuances.
- Da Yin, Hritik Bansal, Masoud Monajatipoor, Liunian Harold Li, and Kai-Wei Chang. 2022. Geomlama: Geo-diverse commonsense probing on multilingual pre-trained language models.

A Appendix

Prompt	Model	Size	US	China	India	Iran	Kenya	Greece	Overall
	LLaMA 2 + English Alpaca	7B 13B 70B	0.28 0.31 0.31	0.50 0.54 0.46	0.39 0.42 0.45	0.24 0.31 0.28	0.31 0.28 0.28	0.34 0.34 0.38	0.34 0.37 0.36
English	Yi + English Alpaca	6B 34B	0.52 0.62	0.43 0.39	0.33 0.42	0.48 0.45	0.50 0.50	0.34 0.44	0.43 0.47
	Uliza + English Alpaca	7B	0.21	0.50	0.39	0.17	0.25	0.31	0.31
	Uliza + {Swahili, English} LoRA	7B	0.45	0.39	0.39	0.34	0.31	0.25	0.36
Hindi LLaMA 2 + Hindi Alpaca		7B 13B 70B	0.28 0.24 0.24	0.46 0.36 0.46	0.39 0.36 0.33	0.34 0.28 0.28	0.25 0.31 0.34	0.41 0.38 0.38	0.36 0.32 0.34
Chinese LLaMA 2 + Chinese Alpaca Yi + Chinese alpaca		7B 13B 70B	0.34 0.38 0.38	0.25 0.32 0.39	0.39 0.39 0.42	0.41 0.48 0.48	0.41 0.47 0.53	0.34 0.38 0.34	0.36 0.40 0.43
		6B 34B	0.38 0.55	0.39 0.54	0.45 0.55	0.34 0.45	0.25 0.44	0.31 0.53	0.36 0.51
Swahili	LLaMA 2 + Swahili Alpaca	7B 13B 70B	0.34 0.34 0.31	0.32 0.29 0.36	0.39 0.39 0.39	0.17 0.24 0.21	0.31 0.34 0.34	0.34 0.34 0.38	0.31 0.33 0.33
	Uliza + Swahili Alpaca	7B	0.31	0.46	0.45	0.28	0.31	0.38	0.37
	Uliza + {Swahili, English} LoRA	7B	0.38	0.32	0.36	0.48	0.41	0.34	0.38
Persian	LLaMA 2 + Persian Alpaca	7B 13B 70B	0.31 0.31 0.28	0.25 0.25 0.36	0.27 0.33 0.33	0.31 0.28 0.14	0.38 0.34 0.25	0.38 0.34 0.38	0.32 0.31 0.29
Greek	Greek LLaMA 2 + Greek Alpaca		0.17 0.21 0.28	0.21 0.21 0.21	0.27 0.30 0.33	0.10 0.17 0.14	0.25 0.28 0.22	0.28 0.28 0.34	0.22 0.24 0.25

Table 6: RQ1, RQ2: Cultural performance scores of various models on the GeoMLaMA benchmark. Values are between 0 and 1, higher is better.

Data collection from native speakers for adapted CAMeL dataset We provided native speakers with a list of words that we procured from different sources on the internet and from large language models as a base collection for each category that they are then asked to verify and correct with more appropriate targets for each category based on their lived experiences.

For the prompts, we follow a similar process, but this time we don't require country specific prompts, only category specific. The final set of prompts is decided by agreement between the authors.

We note that this process has inherent biases for the group of people who perform the tasks, which might implicitly show up in the data in unobserved ways. Also, because two of the categories are about names of people, this may include information about someone's real name, but that would only be so, because it is a common name in some part of their country.

All annotators are demographically located in the USA and are between 25-40 years old. Other than the Hindi annotator who is female, all others identify as male. Also, we note that all annotators are either authors or close friends of authors who did not require any form of compensation.

1000

1001

1002

1003

1004 1005

1006

1007

Category	Size	USA	China	India	Iran	Kenya	Greece
	7	0.03	0.44	0.16	0.34	0.10	0.09
Beverage	13	0.04	0.49	0.16	0.35	0.11	0.09
C C	70	0.04	0.57	0.23	0.38	0.15	0.15
	7	0.09	0.24	0.42	0.08	0.12	0.13
Female Clothing	13	0.11	0.30	0.46	0.09	0.15	0.14
	70	0.10	0.39	0.47	0.11	0.12	0.26
	7	0.05	0.50	0.21	0.28	0.19	0.18
Female Names	13	0.05	0.52	0.32	0.38	0.24	0.33
	70	0.07	0.71	0.33	0.39	0.30	0.37
	7	0.06	0.28	0.06	0.47	0.04	0.23
Food	13	0.07	0.33	0.08	0.49	0.06	0.21
	70	0.18	0.43	0.12	0.53	0.08	0.35
	7	0.10	0.29	0.10	0.34	0.17	0.07
Literature	13	0.12	0.27	0.12	0.39	0.19	0.06
	70	0.16	0.28	0.14	0.65	0.27	0.08
	7	0.09	0.28	0.27	0.39	0.13	0.17
Location	13	0.07	0.35	0.36	0.43	0.18	0.23
	70	0.13	0.39	0.43	0.49	0.19	0.26
	7	0.08	0.62	0.11	0.06	0.17	0.11
Male Clothing	13	0.09	0.68	0.18	0.10	0.19	0.12
	70	0.15	0.73	0.20	0.06	0.16	0.19
	7	0.02	0.53	0.22	0.26	0.30	0.20
Male Names	13	0.04	0.58	0.30	0.33	0.38	0.30
	70	0.04	0.76	0.30	0.35	0.42	0.34
	7	0.16	0.41	0.09	0.11	0.24	0.11
Religion	13	0.15	0.51	0.14	0.14	0.23	0.14
	70	0.17	0.50	0.18	0.16	0.22	0.19
	7	0.08	0.39	0.18	0.26	0.16	0.14
Overall	13	0.08	0.39	0.18	0.26	0.16	0.14
	70	0.08	0.39	0.18	0.26	0.16	0.14

Table 7: RQ3: Setting1 Results (Default MCQ setting, single correct country choice provided) from the CAMeL benchmark.

Prompt	Size	USA	China	India	Iran	Kenya	Greece
	7	0.0	0.36	0.17	0.22	0.14	0.12
USA	13	0.0	0.33	0.18	0.22	0.15	0.12
	70	0.0	0.36	0.16	0.23	0.13	0.12
	7	0.1	0.0	0.24	0.27	0.22	0.17
China	13	0.09	0.0	0.25	0.29	0.21	0.16
	70	0.08	0.0	0.23	0.32	0.18	0.18
	7	0.07	0.37	0.0	0.24	0.18	0.14
India	13	0.07	0.33	0.0	0.27	0.19	0.13
	70	0.06	0.35	0.0	0.29	0.17	0.13
	7	0.09	0.40	0.21	0.0	0.16	0.14
Iran	13	0.08	0.37	0.24	0.0	0.17	0.14
	70	0.07	0.40	0.22	0.0	0.15	0.16
	7	0.08	0.37	0.19	0.23	0.0	0.13
Kenya	13	0.08	0.34	0.21	0.25	0.0	0.13
-	70	0.07	0.34	0.20	0.27	0.0	0.13
	7	0.08	0.37	0.18	0.22	0.15	0.0
Greece	13	0.07	0.33	0.20	0.23	0.16	0.0
	70	0.06	0.36	0.18	0.26	0.14	0.0

Table 8: RQ3: Setting2 Results (Distribution of Countries chosen when correct country is not provided) from the CAMeL benchmark

Category	Size	USA	China	India	Iran	Kenya	Greece
	7	0.48	0.05	0.22	0.09	0.34	0.37
Beverage	13	0.47	0.06	0.24	0.1	0.29	0.32
C C	70	0.5	0.05	0.19	0.08	0.25	0.27
	7	0.35	0.12	0.04	0.4	0.27	0.29
Female Clothing	13	0.37	0.12	0.04	0.35	0.25	0.28
	70	0.36	0.1	0.06	0.36	0.3	0.19
	7	0.48	0.07	0.19	0.16	0.21	0.2
Female Names	13	0.46	0.11	0.14	0.15	0.21	0.19
	70	0.44	0.04	0.2	0.17	0.19	0.17
	7	0.38	0.13	0.37	0.04	0.47	0.13
Food	13	0.37	0.12	0.34	0.04	0.45	0.17
	70	0.31	0.11	0.32	0.05	0.43	0.14
	7	0.3	0.14	0.34	0.1	0.22	0.44
Literature	13	0.28	0.14	0.29	0.07	0.22	0.45
	70	0.32	0.14	0.32	0.05	0.19	0.43
	7	0.42	0.13	0.14	0.11	0.28	0.26
Location	13	0.48	0.11	0.12	0.12	0.25	0.24
	70	0.45	0.13	0.13	0.08	0.29	0.24
	7	0.3	0.02	0.3	0.39	0.28	0.32
Male Clothing	13	0.32	0.02	0.24	0.31	0.25	0.25
	70	0.28	0.02	0.23	0.35	0.29	0.22
	7	0.65	0.06	0.19	0.21	0.16	0.18
Male Names	13	0.67	0.09	0.15	0.2	0.14	0.2
	70	0.67	0.02	0.17	0.2	0.17	0.17
	7	0.27	0.08	0.32	0.36	0.2	0.28
Religion	13	0.27	0.04	0.28	0.3	0.18	0.29
	70	0.27	0.06	0.28	0.28	0.17	0.23
	7	0.4	0.09	0.23	0.21	0.27	0.28
Overall	13	0.41	0.09	0.2	0.18	0.25	0.27
	70	0.4	0.08	0.21	0.18	0.25	0.23

Table 9: RQ3: Setting3 Results from the CAMeL benchmark (How many times did Llama choose the single incorrect option ignoring the other correct options. This number should ideally be 0 for everything.)

Category	Llama_Size	USA	China	India	Iran	Kenya	Greece
	7	0.68	0.72	0.52	0.67	0.79	0.75
Overall	13	0.73	0.71	0.72	0.72	0.78	0.62
	70	0.6	0.61	0.76	0.65	0.7	0.69
	7	0.61	0.47	0.34	0.44	0.76	0.74
beverage	13	0.66	0.43	0.53	0.64	0.68	0.58
	70	0.56	0.31	0.67	0.4	0.5	0.63
	7	0.65	0.83	0.37	0.83	0.85	0.77
female_clothing	13	0.68	0.81	0.62	0.83	0.78	0.57
	70	0.6	0.69	0.58	0.81	0.79	0.69
	7	0.92	0.98	0.8	0.78	0.96	0.93
female_names	13	0.98	0.99	0.93	0.94	0.97	0.77
	70	0.89	0.87	0.97	0.82	0.92	0.85
	7	0.52	0.58	0.3	0.34	0.8	0.35
food	13	0.47	0.46	0.63	0.33	0.8	0.22
	70	0.32	0.4	0.76	0.32	0.69	0.28
	7	0.26	0.37	0.23	0.52	0.41	0.6
literature	13	0.4	0.49	0.38	0.39	0.42	0.56
	70	0.21	0.33	0.45	0.2	0.34	0.65
	7	0.8	0.94	0.6	0.66	0.94	0.93
location	13	0.94	0.93	0.83	0.75	0.97	0.74
	70	0.81	0.88	0.76	0.72	0.84	0.81
	7	0.74	0.65	0.59	0.81	0.8	0.81
male_clothing	13	0.78	0.55	0.82	0.83	0.78	0.6
	70	0.58	0.54	0.85	0.86	0.75	0.74
	7	0.95	0.94	0.78	0.85	0.93	0.91
male_names	13	0.99	0.99	0.93	0.88	0.93	0.87
	70	0.94	0.85	0.97	0.85	0.93	0.87
	7	0.69	0.71	0.63	0.76	0.63	0.71
religion	13	0.63	0.65	0.78	0.78	0.67	0.62
	70	0.51	0.55	0.81	0.72	0.53	0.66

Table 10: RQ3: Setting4 Results from the CAMeL benchmark (How many times did Llama choose an option from the incorrect category) (it was given 3 incorrect categories, 1 correct category) - Ideally this should be 0 for everything if llama understands what category we are asking about.

Category	Size	USA	China	India	Iran	Kenya	Greece
	7	0.37	0.73	0.53	0.26	0.44	0.43
Female Clothing	13	0.38	0.85	0.59	0.34	0.49	0.41
	70	0.4	0.83	0.59	0.31	0.44	0.54
	7	0.12	0.85	0.53	0.52	0.65	0.46
Female Names	13	0.14	0.8	0.64	0.59	0.69	0.51
	70	0.14	0.9	0.58	0.61	0.66	0.58
	7	0.51	0.64	0.79	0.36	0.54	0.59
Male Clothing	13	0.48	0.68	0.8	0.45	0.52	0.56
	70	0.51	0.66	0.77	0.42	0.48	0.68
	7	0.21	0.85	0.59	0.55	0.6	0.56
Male Names	13	0.22	0.82	0.61	0.62	0.57	0.56
	70	0.23	0.93	0.57	0.65	0.68	0.65
	7	0.3	0.82	0.61	0.44	0.55	0.51
Overall	13	0.3	0.8	0.66	0.52	0.57	0.51
	70	0.32	0.89	0.63	0.51	0.56	0.61

Table 11: RQ3: Setting5 Results for the CAMeL benchmark(How many times did Llama choose correct culture but incorrect gender?) (2 options were from correct culture but opposite gender, and 2 options were from incorrect culture but correct gender)

Parameter	Value
Random Seed	42
Number of Epochs	1 (for 34B or 70B models),
	3 (for 6B, 7B, 13B models)
Bits and By	tes Config
Load	4 bit
Quantization Type	nf4
DataType	bfloat16
Lora C	Config
Lora Alpha	16
Lora Dropout	0.1
R	64
Bias	none
Training A	rguments
Per Device Train Batch Size	6 (1 A100 80GB GPU)
Gradient Accumulation Steps	2
Learning Rate	3e-4
Max Gradient Norm	0.3
Warmup Ratio	0.03
Learning Rate Scheduler	constant
Optimizer	32bit paged AdamW
Max Sequence Length	2048

Table 12: Hyperparameters used for Instruction tuning of the LLaMA 2 models