# DIWALI - Diversity and Inclusivity aWare cuLture specific Items for India Dataset and Assessment of LLMs for Cultural Text Adaptation in Indian Context

### Anonymous ACL submission

## Abstract

Large language models (LLMs) are used in various applications. However, despite its wide capabilities, it is shown to lack cultural alignment (Ryan et al., 2024; AlKhamissi et al., 2024) and produce biased generations (Naous et al., 2024) due to a lack of cultural knowledge and competence. Evaluation of LLMs for cultural awareness and alignment is particularly challenging due to the lack of proper evaluation metrics and the unavailability of culturally grounded datasets representing the vast complexity of cultures at the regional and sub-regional levels. Existing datasets for culture-specific items (CSIs) focus primarily on concepts at regional levels and contain several inconsistencies regarding the cultural attribution of items. To address this issue, we created a novel CSI dataset for Indian culture, belonging to 17 cultural facets. The dataset<sup>1</sup> comprises  $\sim 8k$  cultural concepts from 36 sub-regions. To measure cultural competence, we evaluate the adaptation of LLMs to cultural text using the created CSIs, LLMbased, and human evaluations. Also, we perform quantitative analysis demonstrating selective sub-regional coverage and surface-level adaptations across all considered LLMs.

# 1 Introduction

005

017

022

034

Large language models (LLMs), despite having vast knowledge and being trained on extensive data, are shown to lack cultural knowledge and competence (Naous et al., 2024; Wang et al., 2024; Masoud et al., 2025). This is primarily attributed to highly biased pre-training data towards Western culture. Cultural text adaptation refers to modifying text from a general or specific culture to align with a particular cultural group's linguistic, social, and conceptual norms. Cultural adaptation has a wide range of applications, ranging from education (Burstein et al., 2007), dialogue, movie subtitling,



Figure 1: Comparison of the number of Culture-Specific Items (CSIs) between CANDLE (Nguyen et al., 2023) and Ours (*DIWALI*), highlighting differences in CSIs for five facets available in CANDLE with the domain specified as country and subject as India.

hospitality, advertisement, cross-cultural communication, and storytelling. However, ensuring effective and connected cultural adaptation is challenging, requiring understanding and knowledge of social norms, subjective historical references, resonating events, and cultural nuances. These adaptations can generally be surface or deep-level adaptations. We define culture at the regional or country levels, focusing on India. With its vast linguistic and cultural diversity, India presents a unique challenge for culture-driven NLP systems. In particular, India houses 28 states and eight union territories<sup>2</sup>, each with multiple languages, dialects, traditions, and socio-cultural norms. Furthermore, cultural diversity exists at the state level and within sub-regions. For example, the state of Assam has communities like Bodo, Rabha, Karbi, and so on. Each with distinct cultural norms, artifacts, and language<sup>3</sup>.

041

043

045

046

047

048

051

054

060

061

062

Despite recent advances, existing LLMs struggle with cultural adaptation. Studies (Naous et al., 2024) have shown that LLMs trained on Englishcentric corpora often fail to generate culturally ap-

<sup>&</sup>lt;sup>1</sup>Our dataset, code, and LLMs generations will be made available upon acceptance of the work.

<sup>&</sup>lt;sup>2</sup>https://knowindia.india.gov.in/states-uts/ <sup>3</sup>https://assam.gov.in/about-us/391

propriate content. In this work, we evaluate open-063 weight LLMs for cultural text adaptation on closed-064 text generation. In particular, we selected existing 065 text corpora from the education domain originating from the USA. We systematically evaluate multiple open-source LLMs for their ability to perform cultural text adaptation. We rely on CSIs to assess the cultural competence of various LLMs for cultural text adaptation tasks. CSIs represent a particular culture's concepts, objects, items, and customs. These can be utilized to measure the cultural competence of language models. Prior works, such as CANDLE (Nguyen et al., 2023), have made foundational efforts to automatically extract culturally specific concepts from large web 077 text for various cultures. In another work, DOSA (Seth et al., 2024) builds community-driven cultural artifacts for a sub-regional culture belonging to India. However, these datasets have a limited representation of cultural concepts in India. Figure 1 shows the comparison with respect to existing concepts from the CANDLE<sup>4</sup> framework. To address this gap, we built a novel, large-scale dataset of Culture Specific Items for the Indian subregion. In particular, our dataset contains a total of  $\sim 8k$ cultural concepts from 17 facets. We evaluate using automatic CSI-based and human evaluation metrics to assess the adaptation quality of various LLMs systematically. Our dataset provides valuable resources to the community for further research in the comprehensive and standardized assessment of cultural text adaptation. 094 095

In summary, our contributions are as follows:

100

101

102

103

105

106

107

108

109

110

- We present DIWALI, a high-quality Culture Specific Items (CSIs) dataset of India, covering a total of 36 subregions<sup>5</sup>. Our dataset contains a total of 8817 concepts from 17 cultural facets.
- We assess 7 open-weight models across 3 LLM families using both automatic and human evaluation metrics for cultural text adaptation.
- We conduct a cultural subregional evaluation and analysis of LLMs in the context of cultural text adaptation.

#### **Culture & Culture Specific Items** 2

Culture is a multifaceted and complex construct. It means different things to different groups of people. It encompasses knowledge, beliefs, morals, values, and customs. Broadly, we can divide culture into material and non-material components. Material components include tangible elements such as food, dress, houses, and ornaments. In contrast, non-material components refer to symbols, ideas, beliefs, norms, values, morality, and attitudes that guide social behavior. For example, the practice of fasting during Navratris<sup>6</sup>. Culture is a complex concept and hard to define. As a result, the community has turned to using demographic and semantic proxies to approximately define culture (Adilazuarda et al., 2024). Demographic proxies typically include geographical region, language, gender, race, religion, education, and ethnicity. In contrast, people associate semantic proxies with emotional expressions, food habits, social and political relations, actions, and naming conventions. In this work, we consider culture as a regional demographic proxy. With the geographical country India as a region and its states and union territories as sub-regions.

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

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

Newmark (2003) defined culture as "... the way of life and its manifestations that are peculiar to a community that uses particular languages as a means of expression." This definition views culture as a collection of objects, customs, and traditions that hold meaning within a specific community. It categorizes cultural manifestations into five collections or facets: (a) ecology, (b) material culture, (c) social culture, (d) organizations, customs, activities, procedures, concepts, and (e) gestures and habits. CSIs represent concepts, objects, items, and customs that hold significant cultural importance and relevance to a particular culture. It may vary across different subregions. Existing CSI datasets, such as CANDLE and DOSA, are limited in scope, covering only a subset of cultural items. Moreover, automated frameworks like CANDLE suffer from high false-positive rates. CANDLE groups cultural knowledge as a facet and concept. Each facet contains multiple concepts that are culturally relevant to a domain and subject. DOSA (Seth et al., 2024) builds an India-specific cultural dataset by gamifying the process of data collection from 19 socio-diverse geographical locations, with a total of 615 social artifacts. Despite these foundational efforts, through an initial manual inspection, we observe that cultural concepts extracted by CANDLE are incorrectly mapped to a facet (see Appendix

<sup>&</sup>lt;sup>4</sup>https://candle.mpi-inf.mpg.de/

<sup>&</sup>lt;sup>5</sup>Represents states and union territories of India

<sup>&</sup>lt;sup>6</sup>a religious observance common in Northern India.

A for more details). For example: *Light*, *Africa*, 161 Nepal under "Rituals." To filter such false pos-162 itives, we further utilize a prompt-based method 163 to check for cultural relevance with respect to the 164 facets in CANDLE. We design a simple prompt 165 template shown in Prompt - "CSIs Relevancy." For 166 example: Is rishikesh a cultural ritual concept of 167 India? Answer in Yes or No. Each filteration is 168 then manually inspected by one of the authors.

# **3** DIWALI Dataset

170

171

172

173

174

175

176

178

179

181

186

187

189

190

191

193

195

196

197

198

199

204

207

DIWALI comprises 17 systematically curated cultural facets that capture the rich diversity of India across 36 sub-regions. The sub-regional cultural items are defined at the level of states and union territories of the country. Sub-sub-regional distinctions are grouped under their broader sub-region to maintain consistency, e.g., cultural items of communities such as Bodo and Rabha, particularly the indigenous community of Assam, are grouped in the Assam sub-region without further subgrouping. The dataset contains 8817 CSIs, with the highest concept from the food facet, followed by dance forms.

## 3.1 Choice of facets

As discussed in earlier sections, CANDLE considers a limited number of facets that may not be fully representative of the Indian culture. To address this, we expanded the facets to 17 well-known cultural constructs of India. Initially, authors discussed the various cultural concepts that are available across the sub-regions and grouped them into: material, social practices, geographical locations, and names.

## 3.1.1 Material Culture

*Clothing, Textiles, Jewellery, Food,* and *Drinks* cover clothing choices, regional weaving traditions, traditional ornaments, food, and culturally specific traditional drinks.

# 3.1.2 Social Practices and Activities

*Festivals, Rituals, Traditions, Dance forms,* and *Traditional games, Religion, Arts* represent common festivals, religious customs, cultural practices, traditional and prominent dance styles, religious activities, and indigenous games.

3.1.3 Geographical locations and Languages

States & Capitals, Places, Languages & Dialects, Architectural styles covers the territorial region of India, such as states, union territories, state capitals,



Figure 2: Distribution of concepts across different cultural facets in (*DIWALI*).

common places, and architectural styles. It also covers the languages and dialects spoken across various sub-regions.

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

239

240

#### 3.1.4 Names

Each sub-region in India has a distinct naming style pertaining to traditional customs and different tastes for the selection of personal names. Hence, we consider common and significant names. Overall, the 17 facets covers significant concepts

for Indian culture. Concept distributions across facets of DIWALI are presented in Figure 2.

# **3.2** Data construction

We initially start by prompting GPT-40 using the simple prompt mentioned in Appendix section under prompt. Largely, the concepts generated by the prompt contains concepts sourced from Wikipedia. We observed that Wikipedia exhibits significant limitations in representing India's vast cultural diversity. Given this, relying solely on Wikipedia results in an incomplete and less reliable dataset. However, we find that the prompt-based method for curating datasets is limited by the scope of the knowledge learned by the LLM. Given the limitations of normal prompt-based curation, we decided to collect cultural concepts through GPT-40-plus web searches using the detailed prompt in Appendix section under prompt title "Curating India-specific CSIs". To enhance coverage and authenticity, we broaden our data collection scope by searching for official cultural and tourism websites for each sub-region<sup>7</sup>. These sources provide more authentic and verified documentation, ensuring a comprehensive and representative dataset.

<sup>&</sup>lt;sup>7</sup>Here we consider sub-region as the States and the Union territories of India.

241

245

246

247

249

250

255

260

261

262

263

264

270

276

277

278

281

3.3 Quality Check

242One of the authors manually verified the DIWALI243dataset comprising 8817 culture-specific items244(CSIs). The verification steps are as follows:

# 3.3.1 Link validation

We verify the source link of all CSIs. For example, the concept *Kalamkari*, and its corresponding valid source URL: https://en.wikipedia.org/wiki/Kalamkari). Items with invalid source link are excluded from further processing.

# 3.3.2 Concept validation

The same facets were then cross-checked against at least one additional reliable source typically a Government of India portal such as tourism website, archaeological Survey of India. In case of absence of government sources, we relied on Wikipedia. Concepts with incorrect details are removed in this step.

# 4 Cultural Adaptation

# 4.1 **Problem Definition**

We define **cultural text adaptation** as transforming a text x, originally grounded in a source cultural context  $C_s$ , into an adapted text x' that aligns with the target cultural context  $C_t$ . The objective is to preserve the original semantic intent of x while ensuring x' is culturally relevant and emotionally resonant for an audience with  $C_t$ . In this work, we focus on American (source) to Indian (target) cultural contexts and evaluate LLMs' ability for culturally relevant adaptations.

## 4.2 Evaluation Dataset

We use GSM8k (Cobbe et al., 2021) and its multilingual variant MGSM (Shi et al., 2022) as our evaluation datasets. These benchmarks consist of grade school-level math word problems to evaluate the reasoning capabilities of LLMs. We select these datasets for two primary reasons. First, the educational domain is cognitively engaging and often reflects socio-cultural values, making it wellsuited for evaluating cultural adaptation. Second, the problems usually contain culturally grounded entities such as names, food items, locations, and units, which can vary in relevance across cultures. We restrict our analysis to the test set of GSM8k and MGSM.

## 4.3 Models

We evaluate the models from 3 large language model family with model parameters ranging from 1B to 9B: (a) Llama 2 7B Chat (Touvron et al., 2023), (b) Llama 3.1 8B Instruct, (c) Llama 3.2 1B Instruct (Grattafiori et al., 2024), (d) Llama 3.2 3B Instruct (Grattafiori et al., 2024), (e) Mistral 7B Instruct (Jiang et al., 2023), (f) Gemma 2 2B Instruct, and (g) Gemma 2 9B Instruct (Team et al., 2024). We design prompts for cultural text adaptation in both English and Bengali. Specifically, we developed a Assistant-style instruction prompt asking an AI assistant to adapt a given text considering cultural relevance, tone and intent, and cultural sensitivity. Due to the high sensitivity of LLMs towards system instructions and prompts, we provide specific guidelines to follow. For easy processing, we ensure that the adapted text, along with replaced concepts, is in a structured format. Full prompts are presented in the Appendix section. We use the same prompt for all LLMs to maintain consistency. Details on decoding strategies and inference setup are presented in Appendix C.

287

288

290

291

293

294

295

296

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

321

323

324

325

326

327

328

329

330

331

# 4.4 Evaluation Metrics

# 4.4.1 CSI Adaptation Score

To quantify the extent of cultural adaptation in generated text, we propose an *Adaptation Score* that evaluates the validity of concept replacements. Given an input text with a set of replaced cultural concepts:  $R = \{w_1, w_2, \ldots, w_N\}$ , the adaptation score is defined as the fraction of replacements that match a pre-specified set of target cultural concepts. Our matching procedure is performed in two ways: strict matching followed by fuzzy matching. Formally, for each replaced concept  $w \in R$ , we define an indicator function I(w) as:

$$I(w) = \begin{cases} 1, & \text{if } w \text{ is successfully matched,} \\ 0, & \text{otherwise,} \end{cases}$$

We follow two matching strategy:

- 1. Strict Match: We first normalize the concept w (e.g., by lowercasing and removing punctuation) and check for an exact match in the target cultural concept set C. That is, if normalize $(w) \in C$ , then I(w) = 1.
- 2. Fuzzy Match: If a strict match is not found, we perform fuzzy matching using a tokenbased similarity measure<sup>8</sup>. Let  $c^*$  be the

<sup>&</sup>lt;sup>8</sup>We use token\_sort\_ratio metric from the FuzzyWuzzy library with default settings.

334

338

347

351

best fuzzy match from C for w such that similarity(normalize(w),  $c^*$ )  $\geq \tau$ . Here  $\tau$  is a predetermined threshold (e.g.,  $\tau = 80$ in our experiments). If this condition is met, we set I(w) = 1.

Based on the above matching strategy, the Adaptation Score for a sentence is given by: Adaptation score =  $\frac{1}{N} \sum_{w \in R} I(w)$ 

For example: Suppose the LLM adapts three 340 non-Indian concepts: {Muffins, Christmas, Beer} 341 into {Paratha 🗸 , Diwali 🗸 , Ginseng 🗡 }. Assuming two of these replacements are valid matches according to our dataset. The adaptation score is computed as follows:  $\frac{1+1+0}{3} = \frac{2}{3}$ .

# 4.4.2 LLM as Judge

Motivated by recent progress in utilizing large language models to evaluate text generations (Zheng et al., 2023). We conduct our evaluation using two state-of-the-art LLMs: Llama-3.1-8B-Instruct and Mistral 7B Instruct-v0.3. Both models are employed to assess the quality of culturally adapted text based on a detailed prompt and scoring criteria. In particular, we prompt LLMs to score adapted text in three dimensions: Cultural Relevance, Language Fluency, and Mathematically Integrity. We follow a Likert scale of 0 (very poor) to 5 (perfect).

# 4.4.3 Human Evaluation

For human evaluation, we randomly sample 50 359 adapted generations for each of 7 models, totaling 350 instances. Unlike LLM-based automatic evalu-361 ation, performing human evaluation captures subjective cultural nuances that LLMs may overlook. We assess the cultural adaptation of various models along Cultural Relevance dimension, defined in Table 8 Under Appendix D. Each annotator is asked to rate the cultural relevance on a 6-point Likert scale (0 = very poor, 5 = perfect) while assuming "Country: India" as a culture proxy. To ensure consistency, we provide annotators with predefined 370 score descriptions as evaluation instructions mentioned in 10 under Appendix E. We consider two annotators from diverse socio-demographic backgrounds, specifically from the states of Chhattisgarh and West Bengal of India. All evaluators are 375 native Indian residents with 20+ years of residency 377 in their respective sub-regions. The demographic details of the evaluators are presented in Table 11 under Appendix F. Due to resource constraints, we only perform human evaluation for the GSM dataset with an English prompt. 381

#### **Results and Discussions** 5

This section describes results and observations for the considered LLMs on CSIs based on average adaptation score, LLM-as-Judge, and Human Evaluation. Table 1 shows the Average Adaptation and LLM-as-Judge scores.

383

384

385

386

387

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

#### **CANDLE vs. DIWALI: Average** 5.1 **Adaptation Scores**

Table 1 presents a detailed comparison of the average adaptation scores (AAS) obtained using the CANDLE dataset versus our dataset. The table is structured to cover different language prompts (English and Bengali) and GSM and MGSM datasets, with performance metrics provided for both fuzzy and exact matches. To prepare our CSIs for use with Bengali prompts, we performed a zero-shot translation from English into Bengali using GPTo3. Specifically, we prompt each English CSI label: "Translate this term into Bengali without providing additional context or domain knowledge." The result is a direct, one-to-one mapping of our English CSIs into their Bengali equivalents, ready to be matched against Bengali prompts.

For the GSM dataset with an English prompt, our dataset significantly outperforms CANDLE. For example, the Llama-2-7b-chat-hf model achieves a fuzzy match score of 0.615 and an exact match score of 0.855 with our dataset, compared to very low scores of 0.050 and 0.028, respectively, with CANDLE. This clear performance gap is consistently observed across multiple models, indicating that our dataset is more sensitive in capturing the adaptation capabilities of the models.

For the GSM dataset with a Bengali prompt, the differences between the two evaluations remain evident. Although some models (e.g., Llama-2-7b-chat-hf and Llama-3.2-1B-Instruct) have zero scores under both methods, other models exhibit a clear advantage when evaluated with our approach. For example, Llama-3.1-8B-Instruct achieves a fuzzy match score of 0.175 and an exact match score of 0.087 using CANDLE, while our evaluation yields higher scores of 0.403 and 0.131, respectively. Similarly, Mistral-7B-Instruct-v0.3 Gemma-2-9B-Instruct and show notable improvements with our method.

Dataset: GSM									
		Prompt:	English			Prompt: Bengali			
	AAS (C/	ANDLE)	AAS (D	IWALI)	AAS (C/	ANDLE)	AAS (DIWALI)		
Model	Fuzzy Match	Exact Match	Fuzzy Match	Exact Match	Fuzzy Match	Exact Match	Fuzzy Match	Exact Match	
Llama-2-7b-chat-hf	0.050	0.028	0.615	0.855	0.000	0.000	0.000	0.000	
Llama-3.1-8B-Instruct	0.040	0.040	0.400	0.605	0.175	0.087	0.403	0.131	
Llama-3.2-1B-Instruct	0.083	0.058	0.489	0.933	0.000	0.000	0.000	0.000	
Llama-3.2-3B-Instruct	0.051	0.028	0.253	0.672	0.210	0.068	0.337	0.195	
Mistral-7B-Instruct-v0.3	0.106	0.056	0.445	0.563	0.286	0.143	0.531	0.102	
Gemma-2-2B-Instruct	0.067	0.049	0.393	0.635	0.248	0.155	0.456	0.075	
Gemma-2-9B-Instruct	0.060	0.039	0.642	0.479	0.114	0.084	0.365	0.127	
			Dataset	: MGSM					
	1	Prompt:	English			Prompt	Bengali		
	AAS (CA	ANDLE)	AAS (D	IWALI)	AAS (CA	ANDLE)	AAS (DIWALI)		
Model	Fuzzy Match	Exact Match	Fuzzy Match	Exact Match	Fuzzy Match	Exact Match	Fuzzy Match	Exact Match	
Llama-2-7b-chat-hf	0.000	0.000	1.000	1.000	0.051	0.122	0.153	0.612	
Llama-3.1-8B-Instruct	0.040	0.016	0.160	0.494	0.115	0.038	0.231	0.230	
Llama-3.2-1B-Instruct	0.016	0.040	0.200	0.488	0.000	0.000	0.000	0.000	
Llama-3.2-3B-Instruct	0.124	0.077	0.316	0.285	0.094	0.073	0.204	0.058	
Mistral-7B-Instruct-v0.3	0.065	0.024	0.192	0.608	0.254	0.217	0.430	0.094	
Gemma-2-2B-Instruct	0.052	0.013	0.169	0.416	0.118	0.118	0.222	0.033	
Gemma-2-9B-Instruct	0.088	0.036	0.428	0.364	0.032	0.040	0.291	0.085	

Table 1: Comparison of adaptation scores across different models

	Llama		. ·	Mictro					
CR	LIama	MC	CR	LF	MC	Model	Α	B	
3 070	3 10/	1 168	3 201	3 310	1 167	Llama-2-7b-chat-hf	2.02	2.18	
3.182	3.390	4.579	3.246	3.289	4.529	Llama-3.1-8B-Instruct	2.42	2.46	
2.707	2.844	4.136	3.207	3.422	4.576	Llama-3.2-1B-Instruct	0.46	0.72	
3.150	3.314	4.624	3.204	3.323	4.483	Mistral-7B-Instruct	1.86	1.98	
3.561	3.780	4.680	3.351	3.357	4.528	Gemma-2-2B-Instruct	1.90	1.82	
3 236	3 405	4.620	3 270	3 340	4.398	Gemma-2-9B-Instruct	2.38	2.34	
	CR 3.070 3.182 2.707 3.150 3.561 3.398 3.236	Llama           CR         LF           3.070         3.194           3.182         3.390           2.070         2.844           3.150         3.314           3.561         3.780           3.398         3.626           3.236         3.405	Llama           CR         LF         MC           3.070         3.194         4.468           3.182         3.390         4.579           2.707         2.844         4.136           3.150         3.314         4.624           3.561         3.780         4.680           3.398         3.626         4.626           3.236         3.405         4.689	Lama LF         MC         CR           3.070         3.194         4.468         3.291           3.182         3.390         4.579         3.246           2.707         2.844         4.136         3.207           3.150         3.314         4.624         3.204           3.561         3.780         4.680         3.351           3.398         3.626         4.626         3.167           3.236         3.405         4.689         3.277	Lama         MC         Mistra           3.070         3.194         4.468         3.291         3.310           3.182         3.390         4.579         3.246         3.289           2.707         2.844         4.136         3.207         3.422           3.150         3.314         4.624         3.204         3.323           3.561         3.780         4.680         3.311         3.351         3.357           3.398         3.626         4.626         3.167         3.256           3.236         3.405         4.689         3.270         3.340	Lama CR         Lama LF         MC         Mistral           3.070         3.194         4.468         3.291         3.310         4.467           3.182         3.390         4.579         3.246         3.289         4.529           2.707         2.844         4.136         3.207         3.422         4.576           3.150         3.314         4.624         3.204         3.323         4.483           3.561         3.780         4.680         3.311         3.57         4.528           3.398         3.626         4.626         3.167         3.256         4.398           3.236         3.405         4.689         3.270         3.340         4.501	Llama CR         Llama LF         MC         Mistral CR         Mistral LF         MOC           3.070         3.194         4.468         3.291         3.310         4.467           3.182         3.390         4.579         3.246         3.289         4.529           2.707         2.844         4.136         3.207         3.422         4.576           3.150         3.314         4.624         3.204         3.323         4.483           3.561         3.780         4.680         3.351         3.357         4.528           3.398         3.626         4.626         3.167         3.256         4.398           3.236         3.405         4.680         3.270         3.340         4.501	Llama CR         LJama LF         MC         Mistral CR         Model         A           3.070         3.194         4.468         3.291         3.310         4.467         Llama-2-7b-chat-hf         2.02           3.182         3.390         4.579         3.246         3.289         4.529         Llama-3.1-8B-Instruct         2.42           2.707         2.844         4.136         3.207         3.422         4.576         Llama-3.2-1B-Instruct         0.46           3.150         3.314         4.624         3.204         3.323         4.483         Mistral-78-Instruct         1.76           3.398         3.626         4.626         3.167         3.256         4.398         Gemma-2-2B-Instruct         1.90           3.236         3.405         4.689         3.270         3.424         501         Gemma-2-9B-Instruct         2.38	Llama CR         Llama LF         MC         Mistral CR         Mistral LF         MC           3.070         3.194         4.468         3.291         3.310         4.467           3.182         3.390         4.579         3.246         3.289         4.529           2.707         2.844         4.136         3.207         3.422         4.576           3.150         3.314         4.624         3.204         3.323         4.483           3.561         3.780         4.680         3.351         3.357         4.528           3.398         3.626         4.626         3.167         3.256         4.398           3.236         3.405         4.689         3.270         3.340         4.501

(a) LLM-as-Judge scores

(b) Human evaluation scores

Table 2: Evaluation scores of LLM as a judge and Human evaluation for cultural text adaptation.

## 5.2 LLM as Judge Evaluation

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

Table 2a reports the average Likert scores assinged by two LLM judges: Llama and Mistral, for Cultural Relevance (CR), Language Fluency (LF), and Mathematical Integrity (MC). The LLM-as-Judge approach allows for quick and surface-level evaluation of the outputs, capturing observable qualities such as the integration of cultural references, fluency of language, and correctness of mathematical content. The prompt used for LLM-as-Judge are presented in Appendix section under "*LLM-as-Judge Judge prompt*".

Overall, the scores indicate that the evaluated models perform consistently well across these dimensions. For e.g. Gemma-2-9B-Instruct and Mistral-7B-Instruct-v0.3 achieve high ratings consistently under both Llama and Mistral as judge models.Although slight variations exist between the two judges: Gemma-2-9B-Instruct have a slightly higher MC score.

# 5.3 Human Evaluation Results

Table 2b presents the average Cultural Relevance scores as rated by human evaluators for each model on a 6-point Likert scale. *Llama-3.1-8-Instruct* achieves a highest cultural relevancy score of 2.44.

To ensure consistency among annotators and the reliability of the scores. We measure the inter-

annotator reliability using Cohen's  $\kappa$  (Landis and Koch, 1977) on the pair of human raters for each model. As shown in Table 12, there is acceptable agreement for models under the Llama and Gemma families. At the same time, evaluation scores for Mistral-7B have low inter-annotator agreement. 456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

#### 6 Analyses

In this section, we perform a detailed analysis to understand the biases, shallow cultural adaptation, naming stereotypes, and replacement of cultural concepts w.r.t 17 cultural facets of our dataset. **Sub-regional coverage:** Despite their strong incontext capabilities and training on a large corpus LLMs, lack pluralistic (Sorensen et al.) alignment.

**Sub-regional coverage:** Despite their strong incontext capabilities and training on a large corpus LLMs, lack pluralistic (Sorensen et al.) alignment. This leads to the propagation of cultural ideas and views from the learned, causing a dominating effect while underrepresenting other cultures. To understand this better, we try to answer the following question: "Do LLMs, when prompted with the task of regional-level text adaptation, fairly represent all the sub-regional concepts?" To analyze this, we map each replaced concept with a sub-region (State or Union Territory) and plot a geographical heatmap for each facet. The heatmap for food in Figure 3. We observe that there is a significant bias in terms of concepts adapted to represent a sub-region. In particular, adapted concepts belong-

Surface vs. Deep level: Culture is a multi-faceted 490 concept that may vary across individuals and can 491 mean multiple things to multiple people, thus 492 making cultural adaptation inherently challenging. 493 Effective cultural adaptation requires being faithful 494 to the source text, yet connecting while also 495 emotionally resonating with the target audience 496 by incorporating deeper cultural nuances of 497 their respective cultures. (Hershcovich et al., 498 2022) defined culture into three axes: Aboutness, 499 Objectives & Values, and Common ground. In 500 particular, Aboutness refers to the contextual relevance of concepts across cultures. For example, the discussion of "cricket match" might be less meaningful for cultures with very little interest 504 in playing cricket. Thus, cultural adaptation is 506 necessary to adapt from one culture to another by connecting relevant events and scenarios. For example, in the festival of "Durga puja" (event), people from Bengal might relate more to "pandal *hopping*" (scenario), thus impacting the aboutness 510 of the text in context with culture. Such kinds 511 of adaptation are at a deeper level and are more 512 emotionally relatable. To understand the levels of 513 whether a surface or deeper level: "We hypothesize 514 that LLMs struggle with cultural adaptation at a deeper level, particularly in connecting events 516 to meaningful scenarios." To analyze such a 517 deeper level of adaptation, we manually samples 518 of generations of the LLMs. We find that all the LLMs fail to connect events and scenarios 520 together, thus leading to failed "aboutness." In 521 Table 3, we show the detailed analysis of the 522 adapted text for Llama-3.1-8B-Instruct, 523 Mistral-7B-Instruct-v0.1, and 525 Gemma-2-7B-Instruct. In particular, we observe that Llama-3.1-8B-Instruct adapts event Tuesday  $\rightarrow$  Diwali (festival). However, the 527 scenario, i.e., sells CDs, remains the same. The 529 scenario of selling CDs does not emotionally resonate or feel relatable during the festival of Diwali. Thus, we observe that the LLMs are able to perform some level of surface-level adaptation but, fail to make deeper-level adaptations. 533

# 7 Related Work

In this section, we describe the various works related to building specific cultural artifacts in India. CANDLE (Nguyen et al., 2023) develops a global CSI collection by extracting cultural commonsense knowledge from web-scale corpora using an end-to-end method. It identifies candidate sentences using techniques such as string matching, named entity recognition, and a set of handcrafted lexico-syntactic rules. These candidates are then classified into cultural facets-such as geography, religion, or occupation-using zero-shot classification methods based on natural language inference. Finally, CANDLE clusters sentences containing similar facts using sentence embeddings and hierarchical agglomerative clustering, allowing it to produce coherent, representative summaries for each cluster. Although cover large scale culture, CAN-DLE's approach results in only 650 concepts for India, and these entries focus on a limited range of facets like names, food, clothing, etc.

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

583

In contrast, DOSA (Seth et al., 2024) employs a participatory research approach, engaging 260 participants across 19 Indian states to compile a dataset of 615 artifacts. It captures a broader range of cultural dimensions including regional rituals, local customs, and specific art forms providing richer insights into India's subregional diversity.

# 8 Usage of CSIs

Culture-specific items are a part of cultural knowledge bases and help elicit cultural knowledge in various applications, such as (a) cultural competence evaluation: existing literature on cultural competence evaluations of LLMs like Culture Bias Score (Naous et al., 2024) relies on culture-specific items to measure the cultural biases of LLMs for text completion tasks. Thus, having a large-scale Indiaspecific dataset with sub-regional categorization will help better evaluate existing LLMs for cultural competence across various downstream tasks, (b) cross-culture translation: existing machine translation systems are "cultural agnostic"; however, recent work (Conia et al., 2024; Yao et al., 2024; Zhang et al., 2024) has proposed a method to perform culture-level translation using implicit knowledge of entities to be translated. Having large-scale culture-specific items will allow the identification of culture-specific entities for better cross-cultural translation, (c) cultural adaptation: cultural adaptation of text in the form of open text generation,



Figure 3: Comparison of heatmaps for Food facets across different subregions.

Model	Adapted Text	Observation
Llama-3.1-8B-Instruct	Rohan sells CDs. He has 8 customers on Diwali. His first 3 customers buy one CD each. His next 2 customers buy 2 CDs each. His last 3 customers don't buy any CDs. How many CDs did Rohan sell on Diwali?	<ul> <li>✓ proper noun</li> <li>✓ event</li> <li>✗ scenario</li> <li>✗ aboutness</li> </ul>
Mistral-7B-Instruct	Ram sells DVDs. He has 8 customers on Tuesdays (Til-lakshadhi). His first 3 customers buy one DVD each. His next 2 customers buy 2 DVDs each. His last 3 customers don't buy any DVDs. How many DVDs did Ram sell on Tuesdays (Til-lakshadhi)?	<ul> <li>✓ proper noun</li> <li>✓ event</li> <li>✗ scenario</li> <li>✗ aboutness</li> </ul>
Gemma-2-9B-Instruct	Babu sells DVDs. He has 8 customers on Tuesday. His first 3 customers buy one DVD each. His next 2 customers buy 2 DVDs each. His last 3 customers don't buy any DVDs. How many DVDs did Babu sell on Tuesday?	<ul> <li>✓ proper noun</li> <li>✗ event</li> <li>✗ scenario</li> <li>✗ aboutness</li> </ul>
Llama-3.1-8b-Instruct	Rohan and Aaradhya are having a contest to see who can eat the most fruit jam-filled parathas, so they unroll as many as they can find. Unfortunately, someone makes a mistake and Aaradhya's was two parathas wide and 24 parathas long while Rohan's was 3 parathas wide and 14 parathas long. If they both ate their entire amount, how many did they eat on average?	<ul> <li>✓ proper noun</li> <li>✓ event</li> <li>✗ scenario</li> <li>✗ aboutness</li> </ul>
Mistral-7B-Instruct	Ramesh and Meera are participating in a competition to determine who can consume the most jalebis, so they unroll as many as they can find. Unfortunately, an error occurs, and Meera's jalebi stack is twice as wide and 24 stacks high, while Ramesh's is 3 stacks wide and 14 stacks high. If they both consume their entire amount, how many did they eat on average?	<ul> <li>✓ proper noun</li> <li>✓ event</li> <li>✗ scenario</li> <li>✗ aboutness</li> </ul>
Gemma-2-9B-Instruct	Anika and Rohan are having a contest to see who can eat the most laddoos, so they unwrap as many as they can find. Unfortunately, someone makes a mistake and Anika's was two laddoos wide and 24 laddoos long while Rohan's was 3 laddoos wide and 14 laddoos long. If they both ate their entire amount, how many did they eat on average?	<ul> <li>✓ proper noun</li> <li>✓ event</li> <li>✗ scenario</li> <li>✗ aboutness</li> </ul>

Table 3: Surface level adaptation shown for original text: (1) *Billy sells DVDs. He has 8 customers on Tuesday. His first 3 customers buy one DVD each. His next 2 customers buy 2 DVDs each. His last 3 customers don't buy any DVDs. How many DVDs did Billy sell on Tuesday?* (2) *Marcell and Beatrice are having a contest to see who can eat the most fruit roll-ups, so they unroll as many as they can find. Unfortunately, someone makes a mistake and Beatrice's was two roll-ups wide and 24 rolls up long while Marcell's was 3 roll-ups wide and 14 roll-ups long. If they both ate their entire amount, how many did they eat on average?* 

such as writing a story (Bhatt and Diaz, 2024) and dialog-based (Singh et al., 2024), has recently been explored. CSIs can enable a better and holistic adaptation.

# 9 Conclusion

584

589

593

In this work, we build a novel Cultural Specific dataset for Indian culture from different subregions, containing a total of 8.8k concepts. We evaluate LLMs using our newly created dataset for the cultural text adaptation task, showing its improved coverage over the existing dataset. Our human evaluation results suggest that LLMs fail to perform cultural adaptation adequately. Furthermore, we analyze surface adaptation, demonstrating LLMs mostly perform shallow-level adaptations. 594

595

598

600

601

602

603

# Limitations

In this work, we introduce a novel CSI dataset, particularly belonging to Indian culture, covering the sub-regional level granularity. However, our study

has certain limitations. First, our evaluations are restricted to widely used large language models 605 that are trained on multilingual and diverse global 606 datasets, rather than models that are specifically pre-trained or finetuned on country-specific linguistic and cultural contexts. We intend to expand our study to country-specific LLMs. Second, our anal-610 ysis on the surface or deeper level of generations 611 is limited to a few examples. The study of "about-612 ness" in cultural contexts necessitates analysis by 613 a large cultural population. Third, we performed a human evaluation by recruiting 4 annotators from 615 3 diverse socio-demographic regions. However, a 616 multifaceted concept such as cultures requires a 617 more diverse evaluation sourced across all possi-618 ble socio-demographic regions to facilitate a more diverse and fair evaluation.

# Ethical Considerations

621

623

625

627

628

634

641

642

644

645

647

652

Culture-specific items are sourced from publicly available sources, including Wikipedia, government, and tourism websites. We do not store any personal information in the DIWALI dataset. The dataset is intended for academic research and cultural analysis with careful consideration to avoid misrepresentation or community-specific practices. Furthermore, we do not collect any personal information from human evaluators.

> All annotators mentioned in Section 4.4.3 are research students, all with at least a bachelor's or master's degree. The annotation was done as part of their research activity they were hired for. The remuneration was covered under their monthly research assistantship.

We only used ChatGPT-o4-mini for assistance purely with the language of the paper, e.g., paraphrasing, spell-checking, or polishing the author's original content, without suggesting new content.

# References

- Muhammad Farid Adilazuarda, Sagnik Mukherjee, Pradhyumna Lavania, Siddhant Shivdutt Singh, Alham Fikri Aji, Jacki O'Neill, Ashutosh Modi, and Monojit Choudhury. 2024. Towards measuring and modeling "culture" in LLMs: A survey. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 15763–15784, Miami, Florida, USA. Association for Computational Linguistics.
  - Badr AlKhamissi, Muhammad ElNokrashy, Mai Alkhamissi, and Mona Diab. 2024. Investigating

cultural alignment of large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12404–12422, Bangkok, Thailand. Association for Computational Linguistics.

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

- Shaily Bhatt and Fernando Diaz. 2024. Extrinsic evaluation of cultural competence in large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 16055–16074, Miami, Florida, USA. Association for Computational Linguistics.
- Jill Burstein, Jane Shore, John Sabatini, Yong-Won Lee, and Matthew Ventura. 2007. The automated text adaptation tool. In *Proceedings of Human Language Technologies: The Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT)*, pages 3–4, Rochester, New York, USA. Association for Computational Linguistics.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *Preprint*, arXiv:2110.14168.
- Simone Conia, Daniel Lee, Min Li, Umar Farooq Minhas, Saloni Potdar, and Yunyao Li. 2024. Towards cross-cultural machine translation with retrievalaugmented generation from multilingual knowledge graphs. In *EMNLP*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The Ilama 3 herd of models. *Preprint*, arXiv:2407.21783.
- 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.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.

801

802

803

804

805

806

807

808

809

767

768

J. Richard Landis and Gary G. Koch. 1977. The measurement of observer agreement for categorical data. *Biometrics*, 33(1):159–174.

710

712

713

715

716

718

720

721

722

723

733

734

735

740

741

742

743

744

745

746

747

749

750

751

752 753

754

755 756

757

758

759

761

- 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.
- Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. 2024. 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.

Peter Newmark. 2003. A textbook of translation.

- Tuan-Phong Nguyen, Simon Razniewski, Aparna Varde, and Gerhard Weikum. 2023. Extracting cultural commonsense knowledge at scale. In Proceedings of the ACM Web Conference 2023, WWW '23, page 1907–1917, New York, NY, USA. Association for Computing Machinery.
- Michael J Ryan, William Held, and Diyi Yang. 2024. Unintended impacts of LLM alignment on global representation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16121–16140, Bangkok, Thailand. Association for Computational Linguistics.
- Agrima Seth, Sanchit Ahuja, Kalika Bali, and Sunayana Sitaram. 2024. DOSA: A dataset of social artifacts from different Indian geographical subcultures. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 5323–5337, Torino, Italia. ELRA and ICCL.
- Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2022. Language models are multilingual chain-of-thought reasoners. *Preprint*, arXiv:2210.03057.
- Pushpdeep Singh, Mayur Patidar, and Lovekesh Vig. 2024. Translating across cultures: LLMs for intralingual cultural adaptation. In *Proceedings of the 28th Conference on Computational Natural Language Learning*, pages 400–418, Miami, FL, USA. Association for Computational Linguistics.
- Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell L Gordon, Niloofar Mireshghallah, Christopher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, and 1 others. Position: A roadmap to pluralistic alignment. In *Forty-first International Conference on Machine Learning*.

- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, and 179 others. 2024. Gemma 2: Improving open language models at a practical size. *Preprint*, arXiv:2408.00118.
- 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, and 49 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.
- Wenxuan Wang, Wenxiang Jiao, Jingyuan Huang, Ruyi Dai, Jen-tse Huang, Zhaopeng Tu, and Michael Lyu. 2024. 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.
- Binwei Yao, Ming Jiang, Tara Bobinac, Diyi Yang, and Junjie Hu. 2024. Benchmarking machine translation with cultural awareness. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 13078–13096, Miami, Florida, USA. Association for Computational Linguistics.
- Zhonghe Zhang, Xiaoyu He, Vivek Iyer, and Alexandra Birch. 2024. Cultural adaptation of menus: A finegrained approach. In *Proceedings of the Ninth Conference on Machine Translation*, pages 1258–1271, Miami, Florida, USA. Association for Computational Linguistics.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.

# Appendix

# A CANDLE False Positives

Cultural concepts extracted by CANDLE are limited and contain false positives for Indian culture. We remove a total of 393 concepts. The detailed statistics of the concepts removed and kept are shown in Table 4.

Facet	Total	Removed	Kept	Sample removed concepts
Clothing	205	107	98	Western, Pakistan, Seersucker, Bohemian, Kimono, Japanese Fireman, Chinese
				Dragon Robe, Poncho
Drink	157	117	40	China, Gewurztraminer, Japan, Thailand, Pakistan, Sri Lanka, Bangladesh, Eng-
				land, Arabica, Robusta, Ginseng, Espresso
Food	176	78	98	Chinese, Continental, Italian, Southeast, American British, Arabica, Robusta,
				Gewurztraminer, Pakistan, Philippines, Burmese, Mauritius, Chinese Food, Sri
				Lankan, Thai, Singapore, Malay, Pakistani, Nepalese
Rituals	183	50	133	Light, Africa, Nepal
Traditions	191	41	150	American, British, Christmas, Africa, Nepal, Jew, Buddhism

Table 4: False positive cultural concepts from CANDLE framework (Country as India)

# **B DIWALI** Details

The facets per sub-region distributions are listed in Table 5. DIWALI contains facet, concept, description, subregion, and source link. Some of the samples from DIWALI are depicted in Table 6. Table 7 lists out the facets considered and their description, along with examples.

Sub-region	Textiles	Dance	Places	Festivals	Food	Games	States	Traditions	Languages	Arts	Jewellery	Religion	Rituals	Clothing	Architectures	Drinks	Names
AN	3	20	15	31	30	4	1	10	20	6	5	7	9	6	10	6	8
AP	8	43	20	37	37	20	1	31	15	10	12	6	10	28	15	10	36
AR	10	20	24	20	20	10	1	15	24	8	20	9	10	58	10	8	11
AS	10	80	15	67	42	20	1	17	15	8	13	7	14	19	19	9	16
BH	7	30	15	65	35	10	1	14	10	8	15	6	19	20	10	7	10
CH	5	13	10	10	50	2	1	2	8	5	10	7	10	10	10	7	10
CT	5	40	15	29	31	12	1	6	15	8	30	9	19	33	6	6	17
DNH	1	13	15	16	27	0	1	7	9	6	4	6	10	9	10	11	8
DL	5	15	10	15	42	10	1	15	7	5	25	8	10	8	10	10	20
GA	4	30	15	84	44	10	1	30	10	11	14	6	14	17	15	11	13
GJ	8	38	15	20	36	11	1	15	12	10	20	9	10	16	13	10	19
HR	7	19	13	24	41	20	1	13	13	9	16	7	10	10	10	10	10
HP	10	43	20	20	31	20	1	8	13	7	20	6	10	15	10	7	7
JK	10	27	27	10	40	11	2	5	10	10	11	6	10	10	10	9	20
JH	7	30	19	23	40	10	1	6	12	9	25	6	10	10	10	5	14
KA	8	41	20	16	39	15	1	13	19	10	15	7	16	11	10	10	25
KL	10	46	15	10	46	12	1	20	37	10	17	8	20	10	20	10	28
LA	5	17	10	12	45	5	1	12	8	4	10	6	10	10	10	6	15
LD	3	13	10	6	38	0	1	5	4	7	10	6	10	5	5	6	13
MP	8	34	15	10	34	11	1	9	10	10	10	7	10	8	10	10	18
MH	10	27	15	9	45	20	1	10	17	10	20	8	12	10	10	10	27
MN	8	24	19	31	47	10	1	10	15	7	9	5	14	10	10	10	46
ML	8	27	10	25	30	8	1	12	13	10	13	6	10	10	8	10	14
MZ	7	12	15	19	37	10	1	20	15	5	5	7	10	9	10	7	10
NL	7	24	10	10	30	9	1	20	20	7	10	7	10	8	10	5	16
OD	8	28	15	60	47	20	1	20	29	10	11	7	10	10	10	10	17
PY	4	5	10	10	37	8	1	10	7	5	10	6	10	5	10	5	19
PB	9	45	15	10	50	20	1	10	10	9	28	6	10	10	10	10	29
RJ	8	38	20	10	41	10	1	12	15	10	13	6	10	10	10	20	12
SK	3	22	10	8	40	5	1	10	12	10	29	7	10	10	10	6	15
TN	18	49	20	12	86	20	1	15	21	6	19	6	15	8	10	10	12
TG	9	45	10	11	36	10	1	8	7	10	40	6	10	8	10	10	16
TR	7	20	15	11	36	15	1	10	10	6	7	4	9	16	10	12	10
UP	6	43	15	47	43	10	1	48	10	7	10	6	10	10	10	12	27
UK	9	38	20	10	36	10	1	20	15	5	19	6	10	8	10	13	10
WB	10	46	10	65	30	14	1	15	15	10	17	7	10	9	10	10	10

Table 5: Distribution of 17 cultural facets (columns) across 36 States/UTs (rows), using abbreviations: AN = Andaman and Nicobar Islands, AP = Andhra Pradesh, AR = Arunachal Pradesh, AS = Assam, BH = Bihar, CH = Chandigarh, CT = Chhattisgarh, DNH = Dadra and Nagar Haveli and Daman and Diu, DL = Delhi, GA = Goa, GJ = Gujarat, HR = Haryana, HP = Himachal Pradesh, JK = Jammu and Kashmir, JH = Jharkhand, KA = Karnataka, KL = Kerala, LA = Ladakh, LD = Lakshadweep, MP = Madhya Pradesh, MH = Maharashtra, MN = Manipur, ML = Meghalaya, MZ = Mizoram, NL = Nagaland, OD = Odisha, PY = Puducherry, PB = Punjab, RJ = Rajasthan, SK = Sikkim, TN = Tamil Nadu, TG = Telangana, TR = Tripura, UP = Uttar Pradesh, UK = Uttarakhand, WB = West Bengal.

## C Decoding Strategies and Inference Setup

**Decoding Settings:** For all experiments, we set the sampling temperature to 0. A temperature of 0 is used to ensure deterministic output generation, which is crucial for reproducibility and consistent evaluation.

811

810

812 813

814

815

819

Facet	Concept	Description	Subregion	Source link
Clothing	Kupaan	Simple cotton wrap worn by Nyishi men, often with bamboo hats.	Arunachal Pradesh	https://www.indiatravel.app/ traditional-dress-of-arunachal-pradesh/
Drinks	Apong	Rice beer brewed by the Mising tribe.	Assam	<pre>https://diversityassam.com/culture/ apong-a-traditional-rice-beer-of-assam/</pre>
Cuisine	Pulihora	Tamarind rice with spices, often prepared during festivals and special occasions.	Andhra Pradesh	https://www.deccanchronicle.com/lifestyle/food-and- recipes/270323/famous-food-of-andhra-pradesh.html
Rituals	Saptapadi	Bengali wedding ritual where couple takes seven steps on betel leaves.	West Bengal	https://sombitdeyphotography.com/blog/ bengali-marriage-rituals
Traditions	Zoti	Traditional practice of offering cucumbers during monsoon.	Goa	https://aratigoa.wordpress.com/2024/06/20/ understanding-goas-ecology-through-rituals-and-festival
States and Capitals	Bengaluru	India's tech hub blending cosmopolitan vibe with Dravidian roots.	Karnataka	https://en.wikipedia.org/wiki/Bangalore
Dance Forms	Jhumair	Folk dance during harvest season, prevalent in North and Western Odisha.	Odisha	https://en.wikipedia.org/wiki/Folk_dance_forms_of_ Odisha
Places	Ross Island	Former British administrative headquarters with ru- ins.	Andaman and Nicobar Islands	https://en.wikipedia.org/wiki/Ross_Island_(Andaman)
Festivals	Losar	Tibetan New Year celebrated in Lahaul and Spiti with traditional rituals.	Himachal Pradesh	https://www.holidify.com/collections/ festivals-in-himachal-pradesh
Religion	Bon	Ancient animistic faith influencing local traditions.	Ladakh	https://www.lehladakhtourism.com/about-ladakh/ ladakh-religion.html
Languages and Dialects	Maram	Maram is a Sino-Tibetan language spoken by the Maram Naga tribe in the Senapati district. The language is integral to the tribe's cultural identity, with rich oral traditions and customary practices.	Manipur	https://en.wikipedia.org/wiki/Maram_language
Arts	Tikuli Art	Tikuli is a traditional art form from Bihar that in- volves the creation of intricate designs on glass, adorned with gold and silver foils. Historically, it was used to make bindis (forehead decorations) for women. Today, Tikuli at Thas evolved to include decorative items and paintings, reflecting the rich cultural heritage of the region.	Bihar	https://www.memeraki.com/blogs/posts/the-beautiful-arts-of- bihar-manjusha-tikuli-madhubani
Architectures	Hawa Mahal	Situated in Jaipur, the Hawa Mahal, or 'Palace of Winds,' was built in 1799 by Maharaja Sawai Pratap Singh. This five-story structure features 953 small windows, called jharokhas, designed to allow royal ladies to observe street festivities without being seen.	Rajasthan	https://en.wikipedia.org/wiki/Hawa_Mahal
Traditionals Games	Phan Sohlaimung	Players stand at a specific distance with a mark between them. Holding a bamboo pole under their right armpits, they grasp it firmly and try to push each other over the mark.	Tripura	https://en.wikipedia.org/wiki/Tripuri_games_and_ sports
Textiles	Himroo Fabric	Himroo is a traditional fabric from Aurangabad, blending silk and cotton to create a luxurious tex- ture. The weaving technique produces intricate pat- terns, often inspired by Persian designs, reflecting the region's historical connections.	Maharashtra	https://textilevaluechain.in/in-depth-analysis/ articles/traditional-textiles/ traditional-textiles-of-maharashtra
Jewellery	Vaddanam	Gold waist belt; heavy ornament; worn during wed- dings and festivals.	Telangana	https://sribhavanijewels.home.blog/2020/07/02/ telangana-traditional-jewellery/
Names	Devya	Means God's gift.	Jammu and Kashmir	https://www.behindthename.com/submit/names/usage/ dogri

#### Table 6: Samples from DIWALI dataset

Facet	Description
Clothing	Traditional and regional attire worn. E.g. Mekhela Chador
Drink	Beverages with cultural and historical significance, including regional specialties and traditional drinks. E.g. <i>Darjeeling Tea</i>
Cuisine (Food)	Dishes and food items. E.g. Dosa
Rituals	Customs followed in religious, social, and life-cycle events. E.g. Chathurthi Vrat
Traditions	Cultural practices, values, and beliefs. E.g. Gaye holud
States & Capitals	Administrative divisions of India. E.g. <i>Gujarat</i>
Dance forms	Traditional and prominent dance styles. E.g. <i>Kuchipudi</i>
Places	Common geographical and locations. E.g. Dal lake
Festivals	Common festivals and celebrations. E.g. <i>Pongal</i>
Religion	Prominent Religious activities E.g. Buddhism
Languages & Dialects	Spoken languages and dialects E.g. Konda
Arts	Painting, Sculpture, and other creative expressions E.g. Terracotta
Architectures	Architectural styles. E.g. Kesariya Stupa
Traditional games	Indigenous games, and recreational activities. E.g. Gilli Danda
Textiles	Regional Weaving Traditions. E.g. Gamsa
Jewellery	Traditional ornaments, and adornments. E.g. Tora
Names	Common and significant names. E.g. Arjun

# Table 7: Facets considered

Moreover, by not specifying any top-p (nucleus sampling) or top-k parameters, we disable stochastic sampling methods in favour of greedy decoding. This further guarantees that the outputs are stable across different runs. We obtained generations using this approach in a single run for the entire dataset.

Inference Setup: Inference is configured to generate a maximum of 2048 new tokens per prompt, ensuring that the models have sufficient capacity to produce complete and coherent responses. Our experiments were conducted on NVIDIA A100 GPUs with 40GB of memory. Depending on the specific experiment, we used either a single GPU or multiple GPUs to accelerate inference.

# **D** Likert Scale

\_ \_

822

823

824 825

826

827

828

829

830 831 The Likert Scale descriptions for Cultural Relevancy is listed in Table 8 and Language and Mathematical fluency are listed in Table 9.

Score	Rubric	Description
0	Very Poor	No adaptation at all, or only direct transliteration / literal translation with no cultural tailoring.
1	Poor	Adaptations are non-sensical or illogical, introducing absurd elements that misrepresent culture.
2	Fair	Only simple replacement of proper nouns; broader context remains unchanged.
3	Moderate	Proper nouns replaced <i>and</i> surrounding context altered to fit typical Indian usage.
4	Good	Multiple entities adapted, showing a coherent and deeper inte- gration with Indian culture.
5	Perfect	Cultural references are deeply integrated and fully resonant; no further adaptation is possible.

Table 8: Descriptions for the 0–5 Likert scale used in evaluating Cultural Relevance.

Score	Rubric	Description (LF)	Description (MF)
0	Very Poor	Uses completely Western terminology and expressions; no Indian linguistic ele- ments appear.	Problem is mathematically nonsensical after adaptation.
1	Poor	Minimal attempt at adaptation; Indian terms used incorrectly or out of place.	Mathematics is unsound or illogical.
2	Fair	Contains a few Indian terms, but they are often mis-used or feel forced.	Basic arithmetic is correct but poorly integrated with cultural context.
3	Moderate	Some Indian terms used correctly, though overall phrasing lacks natural flow.	Mathematically sound, but cultural context could be better aligned.
4	Good	Effectively incorporates Indian English and terminology with a smooth, natural flow.	Well-structured problems with a clear cultural context.
5	Perfect	Seamlessly blends Indian terminology, expressions, and natural language patterns.	Mathematics is flawlessly integrated with cultural elements while re- maining precise and clear.

Table 9: Likert-scale rubric (0-5) for evaluating Language Fluency (LF) and Mathematical Fluency (MF).

# **E** Instructions for Human Evaluation

In this section, we describe the detailed instructions provided to the human evaluators for annotating cultural adaptation scores for various LLMs. The instruction set was designed after multiple brain storming sessions with authors and annotators all belonging from different sub-regions. Details are provided in Table 10.

832

833

834

835

836

Score	Description	Original Text	Incorrect Adaptation	Original Text	Incorrect Adapta- tion
0	No adaptation or complete translitera- tions/translations.	Toulouse has twice as many sheep as Charleston. Charleston has 4 times as many sheep as Seattle. How many sheep do Toulouse, Charleston, and Seattle have to- gether if Seattle has 20 sheep?	Toulouse has twice as many sheep as Charleston. Charleston has 4 times as many sheep as Seattle. Seattle has 20 sheep. Toulouse has $2 \times 4 \times 20$ sheep. How many sheep do Toulouse, Charleston, and Seattle have to- gether?	James decides to run 3 sprints 3 times a week. He runs 60 meters each sprint. How many total me- ters does he run a week?	Ramesh ka 3 sprints 3 times a week 60 meter par 3 baar run karta hai. Isse kya total meter run karta hai?
1	Non-sensical adaptations (breaks logic or value conver- sion).	Darrell and Allen's ages are in the ratio of 7:11. If their total age now is 162, calculate Allen's age 10 years from now.	Ramesh and Rohan's ages are in the ratio of 7:11. If their total age now is 16,200, calculate Rohan's age 10 years from now.	Dan plants 3 rose bushes. Each rose bush has 25 roses. Each rose has 8 thorns. How many thorns are there to- tal?	Ramesh plants 3 gu- lab jamun plants. Each gulab jamun plant has 25 gulab jamuns. Each gulab jamun has 8 thorns. How many thorns are there total?
2	Simple replace- ment of proper nouns only.	Gretchen has 110 coins. There are 30 more gold coins than silver coins. How many gold coins does Gretchen have?	Ramesh has 110 coins. There are 30 more gold coins than silver coins. How many gold coins does Ramesh have?	-	_
3	Proper nouns with minor con- textual change.	Cody eats three times as many cookies as Amir eats. If Amir eats 5 cookies, how many cookies do both of them eat together?	Rohan eats three times as many laddoos as Priya eats. If Priya eats 5 laddoos, how many laddoos do both of them eat together?	_	_
4	Multiple enti- ties changed with deeper connection to Indian culture.	A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?	A traditional Indian silk saree requires 2 bolts of blue silk and half as much white cotton. How many bolts of fab- ric in total are needed?	-	_
5	Fully culturally resonant adap- tation (deemed necessary).	Harry slept 9 hours last night. His friend James slept only 2/3 of what Harry slept. How many more hours did Harry sleep than James?	Rohan slept 9 hours last night. His friend Vijay slept only 2/3 of what Rohan slept. How many more hours did Rohan sleep than Vijay?		_

Table 10: Instruction for Human evaluation

13

# **F Human Annotators**

838

839

In this section, we describe the socio-demographic location of our human evaluators. Details are provided in Table 11.

Evaluator	Location (Sub-Sub-Region/Sub-Region/Country)	YoR <sup>†</sup>	Educational Qualification
А	Hoogly / West Bengal / India	22	Graduate
В	Raipur / Chhattisgarh / India	28	Post-Graduate
С	Contai / West Bengal / India	23	Post-Graduate

Table 11: Demographic information of human evaluators.  $^{\dagger}$ YoR = years of residence in the sub-region.

# 840 G Inter-Annotator Agreement Scores

Table 12 presents the inter-annotator agreement for the various annotated LLMs.

Model	Inter-annotator scores
Llama-2-7b-chat-hf	0.324
Llama-3.1-8B-Instruct	0.345
Llama-3.2-1B-Instruct	0.358
Llama-3.2-3B-Instruct	0.370
Mistral-7B-Instruct	0.200
Gemma-2-2B-Instruct	0.630
Gemma-2-9B-Instruct	0.557

Table 12: Human annotator agreement for each LLMs.

# 842 H Sub-regional Adaptation Coverage

843 The heatmaps for various facets are shown in Figures 5, 6, 7, 9, and 10.

# 844 I Prompts



Figure 4: Comparison of heatmaps for Architecture and Arts.



(b) Dance Forms

Figure 5: Comparison of heatmaps for Clothing, and Dance Forms.



(b) Jewellery

Figure 6: Comparison of heatmaps for Festivals, Jewellery.



Figure 7: Comparison of heatmaps for Language & Dialects, and Places.



(b) Rituals



Figure 9: Comparison of heatmaps for States & Capitals, and Textiles.



Figure 10: Comparison of heatmaps for Traditional games, and Names.



(a) Festivals

#### Figure 11: Comparison of heatmaps for Festivals

#### Simple Prompt

**Template:** Please provide a comprehensive, deduplicated list of 500 culture-specific {facet} items associated with Indian culture. Ensure that each item is unique and pertains specifically to Indian traditions and practices. Present the output in CSV format.

#### Curating India-specific CSIs

**SYS INST.:** Generate a detailed and comprehensive list of traditional and modern dance forms specific to [State/UT], India, covering various cultural, functional, and social contexts. Present the information in CSV format with the following columns:

Dance Form Name: The name of the dance form.

Description: A brief description (maximum 20 words) of the dance form, highlighting its cultural, functional, or symbolic significance.

Reference Link: A reliable source link for further reading or verification.

Ensure the list includes: Dance forms performed in different contexts, including festivals, weddings, rituals, community gatherings, and stage performances. Representation of dance forms for men, women, and children across diverse social, cultural, and economic settings. Detailed coverage of both traditional folk dances and contemporary adaptations. Include dance forms that are specific to various communities, tribes, and cultural groups within [State/UT]. Representative of regional uniqueness and cultural diversity.

Cultural Text Adaptation

**SYS INST.:** You are an AI assistant tasked with adapting a text to suit a specific cultural context in a particular language while maintaining the original meaning and intent. Your goal is to ensure the text feels natural and appropriate for the target audience, considering cultural nuances, values, and sensitivities. When making these adaptations, follow these key steps:

- 1. Cultural Relevance: Adjust any idioms, metaphors, or cultural references that may not resonate with the target audience. Replace them with culturally appropriate alternatives.
- 2. Tone and Intent: Preserve the original emotional tone and message, even when making cultural adjustments.
- 3. Cultural Sensitivity: Be mindful of topics, words, or phrases that might be sensitive or inappropriate in the target culture.

Specific guidelines:

- Replace foreign names with common Indian names (female for female, male for male). Use a diverse set of names.
- Use Indian locations in place of foreign locations.
- Convert all foreign currencies to Indian Rupees (INR) using a clear symbolic or approximate rate (for example, \$1 = 83). Remove any remaining references to foreign currency (USD, etc.).
- Incorporate Indian traditions, festivals, and cultural practices only if it is contextually appropriate and does not distort the original meaning.
- Use regional-specific terminology and expressions without changing the logical sense of the text.
- Replace foreign food items with Indian equivalents only if it makes sense. For example, "muffins" can become "parathas," but do not replace food items that are already commonly used in India or are essential to the text's logic (e.g., do not replace "eggs" if it's about chickens laying eggs).
- Maintain original mathematical operations and numerical values. Do not show any calculations or provide step-by-step solutions.
- Do not transliterate.
- Do not solve the problem or provide the answer.
- · Do not hallucinate or introduce factual errors.
- Ensure the adapted text is coherent and flows naturally in its new cultural context.
- Provide your response as a single-line JSON array without any line breaks or extra whitespace. The response must be valid JSON that can be parsed directly.

For the replaced\_concepts dictionary:

- ONLY include terms that were actually changed to different terms (e.g., "John": "Ramesh").
- DO NOT include terms that remained the same (e.g., do not include "eggs": "eggs" if it was not changed).
- Use the symbol directly in the values, not Unicode escape sequences.
- Include ONLY the meaningful substitutions you made (e.g., "John":"Ramesh", "muffins":"parathas", "\$10":"830").

Example of correct replaced\_concepts:

- Good: {"John":"Ramesh", "muffins":"parathas", "\$10":"830"}
- Bad: {"eggs":"eggs", "John":"Ramesh", "muffins":"parathas", "\$10":"\u20b9830"}

Provide only the adapted text and replaced words of the problem statement in a valid JSON format, like this example:

[{"cultural\_adapted\_text":"Ramesh bought 5 parathas for 830 and gave 2 to his friend. How much money did he spend per paratha?", "replaced\_concepts":{"John":"Ramesh","muffins":"parathas", "\$10":"830"}}]

USER PROMPT: Adapt the following text to Indian culture: {input text}

Prompt in Bengali for Cultural Text Adaptation
SYS INST.: আপনি একজন এআই সহকারী যার কাজ মূল অর্থ ও উদ্দেশ্য বজায় রেখে টেক্সটকে বাংলা ভাষা ও সংস্কৃতির সাথে খাপ খাইয়ে দেওয়া। আপনার লক্ষ্য টেক্সটটিকে বাংলাভাষী শ্রোতাদের জন্য স্বাভাবিক ও প্রাসঙ্গিক করে তোলা, সাংস্কৃতিক সূক্ষ্মতা, মূল্যবোধ ও সংবেদনশীলতা বিবেচনা করে। নিচের ধাপগুলি অনুসরণ করুন:
১. সাংস্কৃতিক প্রাসঙ্গিকতা: বিদেশি প্রবাদ, রূপক বা সাংস্কৃতিক উল্লেখগুলিকে বাংলার সমতুল্য দিয়ে প্রতিস্থাপন করুন
২. ভাষারীতি: সহজ, প্রমিত বাংলা ব্যবহার করুন (সাধু ভাষা) আঞ্চলিক উপভাষা ছাড়া
৩. সাংস্কৃতিক সংবেদনশীলতা: বাংলার সাংস্কৃতিক রীতিনীতি মেনে চলুন
নির্দেশাবলী: - বিদেশি নামগুলিকে সাধারণ বাংলা নামে পরিবর্তন করুন (মহিলা/পুরুষের জন্য আলাদা)
- বিদেশি স্থানগুলিকে বাংলার স্থাননাম ব্যবহার করুন
- মুদ্রাকে ভারতীয় রুপিতে (₹) রূপান্তর করুন (যেমন: \$১ = ₹৮২)
- বাংলার উৎসব/প্রথা যোগ করুন যেখানে প্রাসঙ্গিক
- পশ্চিমা সংখ্যার (123) বদলে বাংলা সংখ্যা ব্যবহার করুন (১২৩)
- গাণিতিক অপারেশন ও সংখ্যাসমূহ অপরিবর্তিত রাখুন
- উত্তর অবশ্যই বাংলা স্ক্রিপ্টে দিতে হবে
- সিঙ্গেল-লাইন JSON অ্যারে ফরম্যাটে উত্তর দিন
উদাহরণ: [{"cultural_adapted_text": "রমেশ ৫টি পরোটা কিনলেন ৮৩০ টাকায় এবং ২টি বন্ধুকে দিলেন। প্রতি পরোটার দাম কত?","replaced_concepts":{"John":"রমেশ","muffins":"পরোটা","\$10":"৮৩০ টাকা"}}]
USER PROMPT: নিচের টেক্সটর্টি বাংলা ভাষা ও সংস্কৃতিতে রূপান্তর করুন এবং JSON ফরম্যাটে উত্তর দিন: {input text}

Figure 12: Bengali Prompt for Cultural Text Adaptation

l

#### LLM-as-Judge

Consider yourself as an AI expert trained to evaluate the cultural adaptation of a given text. **Original Text:** {original\_text}

Adapted Text: {adapted\_text} Target Culture: Indian Rate each criterion on a 0–5 scale and give a concise justification for each score:

#### 1. Cultural Relevance (0–5):

- 0 No adaptation: retains Western concepts with no cultural shift
- 1 Non-sensical adaptation: culturally inappropriate or absurd elements (e.g., "paratha growing on plants")
- 2 Simple proper-noun swap: only names changed to Indian equivalents
- 3 Names plus limited context change: some culturally relevant objects or settings
- 4 Deep contextual change: multiple entities adapted with clear Indian cultural grounding
- 5 Fully resonant: authentic, natural Indian context with no further improvements possible

#### 2. Language Fluency (0-5):

- 0 No Indian linguistic elements
- 1 Very poor: minimal, incorrect Indian terms
- 2 Poor: few Indian terms, often mis-used
- 3 Moderate: correct terms but slightly forced flow
- 4 Good: smooth Indian English with natural terminology
- 5 Perfect: seamless blend of Indian expressions and natural language patterns

#### 3. Mathematical Integrity (0-5):

- 0 Mathematically nonsensical
- 1 Incorrect mathematics
- 2 Basic math correct but weak cultural tie-in
- 3 Sound math; cultural link could be clearer
- 4 Well-structured math with clear cultural context
- 5 Mathematics flawlessly integrated with cultural elements

Format your response *exactly* like this (replace the numbers with your scores):

Cultural Relevance: 5 Explanation: The adapted text is deeply integrated with Indian culture, accurately reflecting significant traditions and practices.

Language Fluency: 5 Explanation: The text uses natural Indian expressions and terminology, making it easy to read and culturally coherent.

Mathematical Integrity: 5 Explanation: The mathematical problem is presented correctly and is well-integrated into the cultural context.

Do not output any extra text or placeholders such as "[score]".

Figure 13: Prompt for LLM-as-judge