Wunchi-Bench: Benchmarking Language Models on Cultural Reasoning with a Focus on Korean Superstition

Anonymous ACL submission

Abstract

As large language models (LLMs) become 001 key advisors in various domains, their cultural sensitivity and reasoning skills are cru-004 cial in multicultural environments. We introduce Nunchi-Bench, a benchmark designed to evaluate LLMs' cultural understanding, with a focus on Korean superstitions. The benchmark consists of 247 questions spanning 31 topics, assessing factual knowledge, culturally appropriate advice, and situational interpretation. 011 We evaluate multilingual LLMs in both Korean and English to analyze their ability to reason 012 about Korean cultural contexts and how language variations affect performance. To system-014 atically assess cultural reasoning, we propose a novel verification strategy with customized scoring metrics that capture the extent to which 018 models recognize cultural nuances and respond appropriately.

Our findings highlight significant challenges in LLMs' cultural reasoning. While models generally recognize factual information, they struggle to apply it in practical scenarios. Furthermore, explicit cultural framing enhances performance more effectively than relying solely on the language of the prompt. To support further research, we publicly release Nunchi-Bench alongside a leaderboard.¹

1 Introduction

027

034

The accessibility of large language models (LLMs) is expanding into everyday use beyond specialized domains such as legal and psychological counseling. However, their uneven cultural knowledge and contextual understanding are significant limitations (Hershcovich et al., 2022), reducing their effectiveness in multicultural settings. Therefore, integrating cultural insights into problem interpretation and advice is crucial for both fairness and functionality.

Knowledge

In Korea, eating seaweed soup on the day of an exam is believed to cause failure.

Trap Question

My roommate's favorite food is seaweed soup. She has been stressed from studying for exams recently. Would it be okay if I *prepare seaweed soup for her on the day of her exam*?

It's thoughtful of you ... However, ... some Koreans avoid eating seaweed soup before important tests or exams because seaweed is slippery ...

That's a very thoughtful idea! ... It's a gesture that is *sure to be appreciated* and could help boost her mood and energy levels before her exam.

Figure 1: Example generations from GPT-4 Turbo and Gemini Pro 1.5 when asked to give advice on the *Trap question*. LMs may have cultural knowledge to solve multiple-choice questions but they can lack the cultural sensitivity needed to navigate complex cultural contexts.

Previous research benchmarking LLMs' multicultural knowledge has typically verified factual information about specific cultures (Myung et al., 2024; Kim et al., 2024a) or explored the models' embedded values (AlKhamissi et al., 2024; Wang et al., 2024). However, in societies where modernity and tradition coexist, perspectives on cultural norms and the extent to which they are followed can vary significantly. Superstitions, in particular, play a profound role in shaping behavioral patterns, even as their adherence differs across individuals and groups. For example, in Korea, there is a superstition that eating chicken wings may lead to infidelity, with differing beliefs about who might be affected (e.g., only women, only men, etc.). Instead of focusing on detailed cultural knowledge or

¹The code and data will be made available upon acceptance.

Superstition	Writing your name in red ink will bring death.
Multiple-Choice Question	What does writing your name in red ink symbolize in Korean culture?
	Choices: A. It will bring death B. It is a sign of good luck C. It is considered a mark of honor D. It represents a new beginning.
Trap Question	My friend really likes the color red. I'm preparing her birthday cake. Would it be okay to use red for the lettering of her name?
Interpretation Question	My fan asked for an autograph, but all I had was a red pen . When I went to sign and write his name, he hesitated and then refused. I wonder why? He didn't seem like he would resell the autograph.

Table 1: Sample questions from Nunchi-Bench, illustrating the three question types: Multiple-Choice, Trap, and Interpretation. The examples highlight the conditions for Trap and Interpretation questions.

the values of the models, our study aims to assess whether LLMs can interpret scenarios and develop strategies that respect diverse cultural values in real conversational settings.

In this paper, we introduce Nunchi-Bench, a benchmark designed to evaluate LLMs' cultural sensitivity and reasoning in the context of Korean superstitions. The benchmark comprises three distinct task types: (1) Multiple-Choice Questions (MCQs) to assess factual knowledge of Korean superstitions. (2) Trap Questions to evaluate the appropriateness of the model's advice in culturally sensitive scenarios. (3) Interpretation Questions to examine whether models can infer cultural meanings from social interactions.

Nunchi-Bench covers 31 topics and includes 31 MCQs, 92 trap questions, and 124 interpretation questions. To facilitate multilingual model evaluation, we provide versions in both Korean and English. Additionally, for trap and interpretation tasks, we offer versions that either explicitly specify or omit references to the Korean cultural context.

Using this benchmark, we evaluate the cultural sensitivity of diverse LLMs capable of processing Korean text, encompassing both private and open-source models. Additionally, we introduce a novel verification strategy for cultural reasoning in LLMs, proposing a scoring metric that assesses how effectively models recognize cultural context and generate responses aligned with specific superstitions.

083

084

087

090

091

093

094

095

097

099

100

101

102

103

104

105

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

In summary, our main findings are: (1) LLMs struggle to apply cultural knowledge in practical scenarios. (2) Cultural contextual cues in the question enhance the models' ability to deliver appropriate responses. (3) Prompt language alone is less effective than explicitly referencing cultural context for generating culturally informed responses. (4) The quality of language-specific training data is crucial.

2 Construction of Nunchi-Bench

2.1 Superstition Collection

We gather superstitions prevalent in Korea from books and news articles. These superstitions are deeply rooted in the cultural influences of East Asia, particularly from China and Japan, and our collection reflects this blend. We include a broad array of superstitions, both traditional and contemporary, without regard for their origins. To assess how well-known these superstitions are, we conduct a fill-in-the-blank quiz with 33 Korean individuals in their twenties. We select 31 out of 35 topics, only those with an accuracy rate of over 50% in the quiz (See Appendix A for details).

2.2 Question Generation

We design tasks to assess language models' understanding of Korean superstitions. These tasks include: (1) *MCQs* that test factual knowledge about Korean superstitions, (2) *Trap Questions* that evaluate whether LMs can provide culturally respectful advice in superstition-related scenarios, and (3) *Interpretation Questions* that assess whether LMs can explain and reason about the potential cultural contexts relevant to a given situation. Table 1 provides a sample set of these questions, showing the same superstition topic in different formats.

Multiple-Choice Question We adapt the fill-inthe-blank questions from Section 2.1 to develop *MCQs* for 31 Korean superstition topics, primarily as a means of assessing the basic cultural knowledge of LLMs before evaluating their performance on the more complex *Trap* and *Interpretation* questions. To ensure that the multiple-choice options are sufficiently challenging and diverse,

181 182

183

184

185

186

187

188

189

190

191

192

193

194

195

197

we utilize the Multicultural Quiz Platform by Chiu et al. (2024), an AI-human collaboration tool for generating culturally relevant *MCQs*.

131

132

133 134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

151

152

153

154

155

156

157

158

160

161

162

164

165

167

168

169

170

172

173

Trap Question evaluates whether LLMs can provide appropriate advice to a user unfamiliar with Korean culture who unknowingly intends to violate or ignore a Korean superstition. In designing these questions, we apply two key conditions: first, the questions must ask for advice using prompts like "Would it be okay to ...?" or "Should I...?" Second, to add complexity, we include traps that explain why the speaker unknowingly feels compelled to act against the superstition, potentially leading the language model to produce an opposite response if it lacks cultural knowledge (e.g., a friend's favorite color being red, which conflicts with the superstition that writing a name in red ink signifies death). To assess the models' ability to navigate multicultural contexts, we create two versions: one where the relatives or friends are explicitly identified as Korean (Specified) and another where no cultural background is specified (Neutral). Eight topics were excluded due to adaptation challenges (see Appendix A).

Interpretation Question are designed to evaluate whether LMs can understand and interpret the cultural nuances behind reactions in specific scenarios. These scenarios involve negative or ambiguous responses from others, whether as a result of a user's actions or not. The questions prompt the models to explore the reasons and meanings behind these reactions. We apply two key conditions: first, the questions must end with prompts like "Why?" or "What could that mean?" Second, we provide reasoning for the user's actions, along with clues to prevent the models from seeking alternative explanations. Like the trap questions, we create versions where the people reacting are either identified as Korean (*Specified*) or unspecified (*Neutral*).

2.3 Quality Check

174To validate the questions, we recruit twelve Korean175participants, each with over ten years of residency176in Korea. Three participants evaluate each question.177Our aim is to ascertain whether the questions are178relevant to Korean superstitions. We directly ask179participants to assess their relevance using three180options: Not related, Related, or I don't understand

what this means. The results show that out of 256 questions, 247 questions are considered relevant by at least two out of three evaluators.

Question Type	Versions	Accepted Questions (Rate %)	Topics Covered
MCQ	Korean English	31 (100%)	31
Trap	Korean+Specified Korean+Neutral English+Specified English+Neutral	92 (93.87%)	23
Interpretation	Korean+Specified Korean+Neutral English+Specified English+Neutral	124 (97.63%)	31

Table 2: *Nunchi-Bench* Question Statistics. The *Versions* column indicates whether questions are written in Korean or English (Korean, English), and whether the scenarios explicitly identify people as Korean (*Specified*) or do not (*Neutral*).

3 Assessing LLMs with Nunchi-Bench

3.1 Experiment Setup

We utilize Nunchi-Bench to assess the cultural sensitivity of six private and six open-source LMs. In selecting these models, we prioritize diversity in their training data. This includes models primarily trained on native Korean data (e.g., HyperClova X), instruction-tuned models that leverage translated Korean data (e.g., KULLM-v3), and multilingual models that are predominantly focused on English and other languages (e.g., Llama-3 8B Instruct), as shown in Table 3.

Туре	Model	Language
	HuperCLOVA V (HCV 002)	Multilingual
	HyperCLOVA-X (HCX 003)	(Korean-Specialized)
	GPT-3.5 Turbo (0125)	Multilingual
Private	Gemini 1.5 Pro-001	Multilingual
	Claude 3 Opus (20240229)	Multilingual
	Claude 3 Sonnet (20240229)	Multilingual
	Mistral Large (2402)	European Languages*
	Qwen 2.5 7B Instruct	Multilingual
	EXAONE 3.0 7.8B Instruct	Korean, English
Open-	Mistral-7B-Instruct-v0.2	English-focused*
source	KULLM-v3	Korean, English
	Llama-3 8B Instruct	Multilingual
	Llama-3.1 8B Instruct	Multilingual

Table 3: Model Selection for Our Experiment. Models marked with an asterisk (*) are not specifically trained on Korean but are included for comparison purposes

We exclude certain open-source Korean models that showed significantly lower performance in our



Figure 2: Model Performance Across Question Types and Language Versions. MCQ scores (left) are shown for English (blue) and Korean (red). Trap (middle) and Interpretation (right) scores are weighted and categorized into Neutral/Specified versions for English and Korean.

		Gemini 1.5 Pro	Claude 3 Opus	Claude 3 Sonnet	Mistral Large	Mistral -7B	Llama-3 8B	Llama-3.1 8B	Qwen 2.5 7B	GPT-3.5 Turbo	Hyper Clova-X	KULLM -v3	EXAON 3.0 7.8B
MCO	English	30	25	15	26	21	22	25	24	22	21	21	19
MCQ	Korean	28	27	9	22	14	11	13	18	19	27	19	19
	English+Neutral	22	16	1	16	1	-1	6	6	9	0	5	14
Turn	Korean+Neutral	48	48	24	8	-11	-6	-1	-1	-10	48	6	11
тар	English+Specified	92	93	64	63	26	47	44	36	35	20	18	21
	Korean+Specified	63	67	45	18	-4	24	-6	20	-6	56	6	18
	English+Neutral	95	75	73	49	23	16	26	34	55	19	20	29
Interpretation	Korean+Neutral	150	148	87	40	-15	6	-11	29	33	105	16	60
	English+Specified	184	184	150	125	60	57	64	68	133	76	50	66
	Korean+Specified	189	182	105	99	1	20	1	43	75	129	2	108

Table 4: Model Scores by Question Type and Language Version. MCQ scores are summed, while Trap and Interpretation scores are weighted. Higher values indicate better performance.

preliminary tests, such as Mi:dm (KT, 2023) and ChatSKKU². For detailed information on the models and the inference methods used, please refer to Appendix B.

3.2 Evaluation Setup

198

199

202

204

207

208

211

212

213

214

215

216

217

For *MCQs*, we calculate accuracy by comparing the model's output with the correct answer. If the model refuses to provide an answer or generates a response in a language other than Korean or English, we mark the response as incorrect.

Evaluating responses to *Trap* and *Interpretation Questions* requires a more nuanced approach. To address this, we develop a specialized scoring system that focuses on the model's cultural sensitivity and its ability to understand specific superstitions.

- 0 points: The response does not mention cultural differences.
- 1 point: The response acknowledges cultural differences but does not directly address the superstition in question.

• 2 points: The response acknowledges cultural differences and accurately relates to the specific superstition.

218

219

221

222

223

224

225

226

227

229

230

231

232

233

234 235

236

237

• -1 point: The response mentions cultural differences but includes incorrect or irrelevant information about the superstition.

This metric is intimately related to the evaluation and verification of rationales generated by LLMs, especially for cultural reasoning focused on the cultural aspects. We employ this metric to evaluate the responses of the models, utilizing GPT-4 Turbo (0409) as the Evaluator. For the details, refer to Appendix C.

3.3 Results

Figure 2 and Table 4 show model performance across question types and language versions. We find that:

Gemini 1.5 Pro and Claude 3 Opus lead Claude 3 Opus and Gemini 1.5 Pro consistently achieved the highest scores across all three question types (MCQ, Trap, and Interpretation), particularly in

²https://huggingface.co/jojo0217/ChatSKKU5.8B



Figure 3: Breakdown of Model Scores in Trap and Interpretation Tasks

English versions.

240

241

242

243

246

247

248

252

260

263

265

267

268

Prompt language impact varies by model Except for HyperClova-X and EXAONE, all models performed better in *English MCQ* than in *Korean MCQ*, indicating a preference for English prompts. The influence of language is also evident in Trap and Interpretation tasks: HyperClova-X excelled in *Korean+Specified* Trap, while both HyperClova-X and EXAONE led in *Korean+Specified* Interpretation. In contrast, all other models performed best in *English+Specified* versions for both tasks. Since HyperClova-X and EXAONE are primarily trained in Korean, this suggests that language-specific training significantly influences model performance, a point further explored in the discussion section.

Cultural cues enhance Trap and Interpretation performance Providing explicit cultural context significantly enhances model performance, with *Korean+Specified* outperforming *Korean+Neutral* and *English+Specified* surpassing *English+Neutral* across most models, except for Llama-3.1 and KULLM-v3. Notably, *English+Specified* exceeds *Korean+Neutral*, suggesting that contextual framing contributes more to reasoning performance than the language of the prompt itself.

269Lower scores in Korean versions relative to English+Neutral Since prompt language provides270glish+Neutral Since prompt language provides271context, English+Neutral contains the least cul-272tural information among the four versions. How-273ever, in Trap and Interpretation tasks, some models274scored lower in the Korean version than in English+Neutral. This is due to receiving -1 scores

from hallucinations, which will be further discussed in the following section.



Figure 4: Score Composition of Models on Trap and Interpretation Questions across Different Versions

3.4 Score Composition Analysis

Figures 3 and 4 illustrate the score composition for *Trap* and *Interpretation* questions across different versions, reinforcing the findings presented in the results. Our analysis reveals the following:

278

279

Culturally specified prompts enhance cultural knowledge retrieval. Expanding on our overall score analysis, we find that cultural specification consistently increases Score 2 counts, allowing models to retrieve and articulate cultural knowledge more effectively. This effect is particularly pronounced in English, where Specified prompts lead to a notable performance improvement, while Korean prompts elicit stronger cultural responses overall.

285

295

312

313

314

315

317

318

319

321

323

327

329

Increased hallucination in Korean versions In

Trap questions, all models-except Claude Sonnet and HyperClova-X-exhibit higher hallucination rates in the Korean versions compared to their English counterparts. This accounts for the lower weighted scores observed in the Korean versions relative to English+Neutral. A similar pattern emerges in Interpretation questions, where hallucination rates in Korean+Neutral exceed those in English+Neutral for all models except Claude Opus and HyperClova-X.

Increased cultural reasoning in Interpretation

questions As shown in Figure 4, the frequency of culturally relevant response attempts (scores other than 0) increases across all versions in Interpretation questions compared to Trap questions. This trend arises because, unlike Trap questions, which assess whether a scenario is problematic, Interpre*tation* questions examine why it is problematic. As a result, models more frequently cite cultural differences as justifications, leading to a higher oc-316 currence of culturally informed responses.

3.5 Error Analysis

Figure 5 presents the number of models failed across topics and question types. Trap questions are the most challenging, followed by Interpretation and MCQ, indicating that models struggle to apply superstition knowledge to cultural scenarios.

For MCQ, the most difficult topic in both versions is ID 16 ("In Korean culture, what are you believed to turn into if you lie down immediately after having a meal?"). While the correct answer is "a cow," most models select "a pig" or "a snake."

For Trap questions, all models fail on ID 25, which is particularly noteworthy given that, in the corresponding MCQ, most models correctly select C ("As poison") for the question: "In Korean cul-



Figure 5: Model Failure Count by Superstition Topics and Question Types. For MCQs, a model fails if it selects an incorrect answer. In Trap and Interpretation questions, failure occurs when no response within the topic scores 2, and the combined count of -1 and 0 scores exceeds that of 1 scores.

ture, how are apples eaten in the evening regarded?" (Options: A. As a healthy snack, B. As a dessert, C. As poison, D. As a traditional medicine). However, when presented with a practical scenario-such as serving apples to Korean relatives at night-models fail to recognize the cultural context.

Similarly, for ID 15 ("Eating Seaweed Soup on Exam Day Brings Failure"), most models correctly identify the superstition in MCQ but struggle in *Trap* questions when required to apply this knowledge to real-world decision-making-such as preparing a meal for someone taking a test (as shown in Table 5).

In *Interpretation* questions, the gap between the Specified and Neutral versions is more pronounced than in Trap questions, with the highest number of model failures occurring in the English+Neutral version. Interestingly, for ID 25, while all models fail in Trap questions regardless of version, several models in the Specified version of the Interpretation question correctly interpret the situation based on the superstition (as shown in Table 6). This sug-

353

333

Question	My Korean roommate's favorite food is seaweed soup. She has been stressed from studying for exams recently. Would it be okay if I prepare seaweed soup for her on the day of her exam? (<i>English+Specified</i>)
Incorrect	Yes, including seaweed soup in your feature on light and nutritious foods for important days like exam days is a great idea (GPT3.5 Turbo, 0-point)
Ideal	It's important to note that some Koreans avoid eating seaweed soup before important tests or exams be- cause seaweed is slippery, and there's a superstition

Table 5: Examples of responses to Question ID 15, with the ideal answer from GPT-4 Turbo (evaluator).

gests that models' cultural reasoning is influenced by both question type and the availability of cultural context.

Trap	Apples are a great choice for dessert, especially if you know that your mother-in-law enjoys fruit and has no allergies to apples (0-point)
Interpretation	In Korean culture, there is a belief that eating apples in the evening or at night can cause indigestion or abdominal discomfort (2-point)

Table 6: Examples of responses to a question in topic ID 25 from Claude 3 Opus in the *English+Specified* version. The Trap response lacks cultural awareness, whereas the Interpretation response incorporates cultural knowledge.

4 Discussion

Are MCQ scores correlated with Trap and Interpretation scores? As shown in Figure 6, English MCQ scores correlate only with the *English+Neutral* Trap questions, while Korean MCQ scores exhibit broader correlations across multiple Trap (*English+Neutral*, *Korean+Neutral*) and Interpretation (*Korean+Neutral*, *Korean+Specified*) versions. No other significant correlations were found.

However, when examined within individual superstition topics, no consistent pattern emerges. Figure 7 illustrates the Spearman correlation between Korean MCQ scores and *Korean+Neutral* Trap questions, revealing fluctuations along the diagonal, where correlations within the same topic vary unpredictably. This inconsistency underscores the limitations of MCQs in assessing cultural reasoning, suggesting that they fail to capture deeper contextual understanding. For all correlation plots and statistics between MCQ and Trap/Interpretation scores, see Appendix D.

What Is the Impact of Korean Language Training on Cultural Reasoning? Figure 8 shows



Figure 6: Statistically significant correlations between MCQ scores and Trap/Interpretation scores across different versions. Pearson and Spearman coefficients are reported for each condition.



Figure 7: Spearman correlation between Korean MCQ scores and *Korean+Neutral* Trap question scores by superstition topic

model performance across three metrics: non-zero score count, weighted sum, and positive score count. Two key trends emerge in *Trap* questions.

First, HyperClova-X consistently outperforms in

359

361

367

371

374

375

the *Korean+Specified* version across all metrics. As a private Korean-focused model, it highlights how high-quality Korean language training enhances both sensitivity and performance in cultural reasoning for Korean prompts.

Second, GPT-3.5, HyperClova-X, KULLM-v3, and EXAONE generate as many or more nonzero responses in *Korean+Specified* than in *English+Specified*, while other models show the opposite trend. Since KULLM-v3 and EXAONE are open-source Korean-focused models, this suggests that language-specific training boosts sensitivity to Korean prompts. However, this does not necessarily improve performance, as seen in the weighted sum and positive score count.

For *Interpretation* questions, the first trend persists, but the second does not. As noted earlier, this likely stems from the nature of *Interpretation* questions, where the tendency to provide culturally relevant responses increases across all versions.

5 Related Work

386

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

429

430

431

432

433

434

435

Research on benchmarking multicultural knowledge in LMs focuses on factual knowledge and embedded values. Studies on factual knowledge assess how well LMs capture culturally specific information. Kim et al. (2024a) introduce CLIcK to evaluate LLMs' understanding of Korean culture, revealing significant gaps, especially in opensource models. Liu et al. (2024) find that multilingual models struggle with proverbs, particularly in cross-cultural and figurative tasks. Myung et al. (2024) assess cultural knowledge across 16 countries with BLEND benchmark, highlighting performance gaps in underrepresented regions.

In contrast, research on embedded values examines biases and cultural alignment. Wang et al. (2024) identify cultural dominance in LLMs, showing a bias toward English-centric norms, even in non-English queries, and propose more diverse pretraining to address this. AlKhamissi et al. (2024) also explore cultural alignment, highlighting Western biases and proposing Anthropological Prompting to improve models' cultural sensitivity.

Studies on the rationales generated by LLMs include verification of the rationales via specific prompts. Vacareanu et al. (2024) propose general principles that a model should follow while reasoning (relevance, mathematical accuracy, logical constituency) to evaluate the model's reasoning chains. Fayyaz et al. (2024) study LLMs' rationales from



Figure 8: Heatmap of Model Performance on Trap and Interpretation Questions. This heatmap compares model performance across three metrics: Non-Zero Score Count (scores of 2, 1, or -1, indicating an attempt at a culturally relevant response), Weighted Sum (aggregated score), and Positive Score Count (scores of 2 or 1). Columns represent each version, with color intensity standardized within each model.

their decision-making process, prompting the models to identify the most important words in the input texts.

6 Conclusion

This study introduced Nunchi-Bench, a benchmark for evaluating LLMs' cultural sensitivity and reasoning, with a focus on Korean superstitions. Our findings reveal significant disparities in how LLMs handle culturally nuanced questions, influenced by question type, prompt language, and the presence of explicit cultural context.

To foster further research, we publicly release Nunchi-Bench and a leaderboard, encouraging ongoing improvements in LLMs' cultural understanding. Future work should extend this benchmark to diverse cultural contexts, ensuring AI systems are not only multilingual but also culturally adaptive.

451

453 454

455

456 457

458

459

460

461

462

463

464

465

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

Limitations

While our study provides valuable insights into the cultural sensitivity of LLMs within Korean contexts, several limitations must be acknowledged.

Cultural Scope Nunchi-Bench is specifically designed to assess cultural reasoning in the context of Korean superstitions. While this focus enables a deep and nuanced evaluation of LLMs in this domain, it limits the generalizability of our findings to other cultural settings. Future research should extend the benchmark to additional cultural traditions and belief systems to enable a more comprehensive assessment of LLMs' cultural adaptability.

Model Specificity Our evaluation includes a selection of contemporary private and open-source multilingual models. However, given the rapid evolution of LLMs, our findings may not generalize to future models that incorporate different training paradigms, larger datasets, or novel architectures. Continuous benchmarking and updates will be necessary to track improvements in cultural reasoning capabilities.

Evaluation Methodology The scoring system for Trap and Interpretation questions relies on a verification strategy using GPT-4 Turbo as the evaluator. While efforts were made to refine this evaluation process through multiple iterations, potential biases in the evaluator model and the scoring framework may influence the results.

Ethics Statement

In our research, we committed to strict ethical standards to ensure inclusivity and fairness. Evaluating language models on their capability to process culturally specific content raises sensitive cultural issues. To mitigate ethical concerns, we meticulously designed the benchmark to prevent the reinforcement of stereotypes and to prompt models to exhibit a nuanced comprehension of cultural variations, rather than just superficial recognition.

By releasing Nunchi-Bench and its leaderboard to the public, we promote transparency and encourage the broader AI research community to participate in developing culturally aware AI technologies responsibly. This open access strategy enhances peer review and fosters the integration of ethical practices by providing resources that can help audit and refine AI systems according to culturally sensitive standards.

References

- 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.
- AI Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. *In Anthropic Model Card*, 1.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.
- Yu Ying Chiu, Liwei Jiang, Maria Antoniak, Chan Young Park, Shuyue Stella Li, Mehar Bhatia, Sahithya Ravi, Yulia Tsvetkov, Vered Shwartz, and Yejin Choi. 2024. Culturalteaming: Ai-assisted interactive red-teaming for challenging llms' (lack of) multicultural knowledge.
- Mohsen Fayyaz, Fan Yin, Jiao Sun, and Nanyun Peng. 2024. Evaluating human alignment and model faith-fulness of llm rationale. *ArXiv*, abs/2407.00219.
- Gemini-Team-Google. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context.
- 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.
- Eunsu Kim, Juyoung Suk, Philhoon Oh, Haneul Yoo, James Thorne, and Alice Oh. 2024a. Click: A benchmark dataset of cultural and linguistic intelligence in korean. *arXiv preprint arXiv:2403.06412*.
- Jeongwook Kim, Taemin Lee, Yoonna Jang, Hyeonseok 554 Moon, Suhyune Son, Seungyoon Lee, and Dongjun 555

503 504

505

506

507

508

509

514

515

516

517

518

519

520

521

522

523

524

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

- 556 557
- 559

- 562
- 563 564
- 565 566 567 570
- 571 572
- 574
- 575
- 576

- 583

584

- 585 586
- 587

601

606

Kim. 2024b. Kullm3: Korea university large language model 3. https://github.com/nlpai-lab/ kullm.

- KT. 2023. Mi:dm: Kt bilingual (korean, english) generative pre-trained transformer. https://genielabs. ai.
- LG-AI-Research. 2024. Exaone 3.0 7.8b instruction tuned language model. arXiv preprint arXiv:2408.03541.
- Chen Liu, Fajri Koto, Timothy Baldwin, and Iryna Gurevych. 2024. Are multilingual LLMs culturallydiverse reasoners? an investigation into multicultural proverbs and sayings. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 2016–2039, Mexico City, Mexico. Association for Computational Linguistics.

Llama-Team. 2024. The llama 3 herd of models.

- Junho Myung, Nayeon Lee, Yi Zhou, Jiho Jin, Rifki Afina Putri, Dimosthenis Antypas, Hsuvas Borkakoty, Eunsu Kim, Carla Perez-Almendros, Abinew Ali Ayele, et al. 2024. Blend: A benchmark for llms on everyday knowledge in diverse cultures and languages. arXiv preprint arXiv:2406.09948.
- Qwen-Team. 2024. Qwen2.5: A party of foundation models.
 - Naver HyperCLOVA X Team. 2024. Hyperclova x technical report.
 - Robert Vacareanu, Anurag Pratik, Evangelia Spiliopoulou, Zheng Qi, Giovanni Paolini, Neha Anna John, Jie Ma, Yassine Benajiba, and Miguel Ballesteros. 2024. General purpose verification for chain of thought prompting. ArXiv, abs/2405.00204.
 - Yuhang Wang, Yanxu Zhu, Chao Kong, Shuyu Wei, Xiaoyuan Yi, Xing Xie, and Jitao Sang. 2024. CDEval: A benchmark for measuring the cultural dimensions of large language models. In Proceedings of the 2nd Workshop on Cross-Cultural Considerations in *NLP*, pages 1–16, Bangkok, Thailand. Association for Computational Linguistics.

Α **Collection Details**

Table 7 presents the fill-in-the-blank questions on Korean superstitions along with their correct answer rates. Only questions with a correct answer rate exceeding 50% are included in the final benchmark. When multiple correct answers were possible, any of the valid options were accepted. Figure 9 illustrates the template used for this purpose. For Trap Ouestions, topic IDs 9, 12, 13, 16, 21, 26, 31, and 32 were excluded due to the difficulty of adapting those topics to the question format.

id	Fill-in-the-blank Question	Correct Answer	Correct Answer Rate (N=33)	
0	숫자 _ 는 불길하다. The number _ is unlucky.	4	100	
1	를 틀고 자면 죽는다. Sleeping with on can cause death.	선풍기 fan	100	
2	밤에 피리나 휘파람을 불면 _이 나온다. Whistling or playing a flute at night brings out	뱀/귀신 snake/ghost	90.91	
3	연인에게를 선물하면 도망간다. Giving to your lover will make them leave.	신발/구두 shoes	90.91	
4	색으로 이름을 쓰면 죽는다. If you write a name in color, the person will die.	빨간/붉은 red	100	
5	연인과 길을 걸으면 헤어진다. Walking on with your lover causes a breakup.	덕수궁 돌담 Deoksugung Path	57.58	
6	국에가 나오면 돈이 생긴다. Dream of a, and you'll receive money.	돼지 pig	93.94	
7	닭 를 먹으면 바람난다. Eating a chicken makes a person flirtatious.	날개 wing	66.67	
8	드 프 기 · · · · · · · · · · · · · · · · · ·	나비/나방 butterfly/moth	12.12*	
9**	아이 위를 넘어다니면, _가 안 큰다.	7] 970w	66.67	
10	는 행운을 가져온다. (힌트: 새) brings good luck (Hint: bird)	까치 magnie	60.61	
11	Shings good tack (Hint: Shin) 소리는 불운을 가져온다. (힌트: 새) The sound of brings had luck (Hint: bird)	까마귀	93.94	
12**	별면 복이 달아난다.	다리	100	
13**	날으면 복이 나간다. (힌트: 실내) Stanning on ruins your luck (Hint: indoore)	문지방/문턱 threshold	81.82	
14	에 순발톱을 깎으면 안된다. 왜냐하면 _가 먹고 사람이 될 수 있기 때문이다. You shouldn't cut your nails at _, because, can set them and turn into you	밤, 쥐 night, rat	72.73	
15	시험날에을 먹으면 시험에 떨어진다.	미역국 seaweed soun	100	
16**	법 먹고 바로 누우면 _가 된다.	소 cow	81.82	
17	법 먹을 때 상의에 앉아서 먹으면 안된다.	모서리 corner	66.67	
18	아들을 낳은 여성의을 입으면 아들을 낳는다. Wear a son-bearing woman'sto have a son	속옷 underwear	18.18*	
19	중에는 장례식에 가지 않는다. You should not attend funerals during	임신 pregnancy	54.55	
20	잡귀를 쫓기 위해서을 뿌린다.	소금/팥 salt/red beans	90.91	
21**	을 먹으면 노래를 잘하게 된다.	날달걀	81.82	
22	· · · · · · · · · · · · · · · · · · ·	태명 pra hirth name	54.55	
23	때 아기가 잡는 물건이 아이의 장래를 나타낸다.	돌잡이 Daliahi	100	
24	지, the them a baby graps shows their jutate.	수직으로 Stishing warting//w	87.88	
25	에 사과를 먹으면 독사과가 된다. (힌트: 시간)	방	81.82	
26**	An apple eaten atis poisonous. (Hint: time) 산성비를 맞으면가 된다.	nigni 대머리	96.97	
27	· J you ger caught in acta rain, you will ger	도깨비	21.21*	
28	Pray to to find lost items.	goblin क्	69.7	
29	If you get a tot of you it tive a tong life. 부조금은 수여야 하고, 단위로 내서는 안 된다. You must give condolence money in numbers and not in units of	eurse words 홀수, 천원 odd, 1000 won	54.55	
30	·····································	흰 국화 white chrysanthemum	75.76	
31**	나이가로 끝날 때를 조심해야 한다. Be careful when your age ends in	9	69.7	
32**	웃다가 웃으면 엉덩이에 난다. Laugh after crying, you'll get on your butt	털/뿔 hair/horns	96.97	
33	머리는 _쪽으로 놓고 자서는 안된다. Don't sleep with your head facing	북/문 north/door	87.88	
34	를 구우면 생기는 물은 보약이다. Water from grilled is medicine.	버섯 mushrooms	24.24*	

Table 7: Fill-in-the-Blank Quiz on Korean Superstitions (Correct Answer Rate in %) Questions marked with an asterisk (*) were excluded from Nunchi-Bench, while topics marked with two asterisks (**) were not included in Trap questions due to adaptation difficulties. Hints were provided only for questions considered overly ambiguous.

제시된 문장 속 빈칸 채워넣기
아래 한국 미신과 관련된 문장에 대해 빈칸을 채워 넣어주세요. Please fill in the blanks for the following sentences related to Korean superstitions. ♥ 조사 '윤/를', '이/가'는 힌트가 아닙니다. (e.g. '클로 끝나니 모음으로 끝나는 단어인가? (X)) ♥ 모르는 항목에 대해서는 .을 찍고 넘어가주시길 바랍니다. ♥ 여러 개의 정답이 가능할 수 있습니다. 둘 이상일 감우에는 /로 구분해서 넣어주세요 The particles '윤/를' and '이/가' are not clues. If you're unsure of the answer, simply mark it with a dot (.) and move on. There may be multiple correct answers. If so, separate each answer with a slash (/).
[1/36] 숫자 _ 는 불길하다 *
The number _ is unlucky.
내답변
[2/36]를 틀고 자연 죽는다 * Sleeping with on can cause death. 내 답변
[3/36] 밤에 피리나 휘파람을 불면 _이 나온다.* Whistling or playing a flute at night brings out 내 답변
[4/36] 연인에게를 선물하면 도망간다 * Walking on with your lover causes a breakup. 내 답변
[5/36]색으로 이름을 쓰면 죽는다 * If you write a name in color, the person will die. 내 답변

Figure 9: Template used for the fill-in-the-blank quiz on Korean superstitions

B Experiment Details

609

610

612

613

614

616

617

618

619

622

623

625

626

The private models we evaluate are as follows:

- HyperCLOVA-X (HCX 003): The leading model in the HyperCLOVA-X series, created by the Naver (Team, 2024).
- Mistral Large (2402): The model developed by Mistral AI³
- GPT-3.5 Turbo (0125): The latest model in the GPT-3.5 series from OpenAI (Brown et al., 2020).
- Gemini 1.5 Pro-001: A multilingual language model from Google (Gemini-Team-Google, 2024),
- Claude 3 Opus (20240229): The largest model from Anthropic (Anthropic, 2024),
- Claude 3 Sonnet (20240229): The second-largest model from (Anthropic, 2024).
- The open-source models we evaluate are as follows:

• Qwen 2.5 7B Instruct: One of the latest models in the Qwen large language series developed by the Qwen Team at Alibaba Group (Qwen-Team, 2024).

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

- EXAONE 3.0 7.8B Instruct: the latest bilingual model in the EXAONE series developed by LG AI Research (LG-AI-Research, 2024).
- Mistral-7B-Instruct-v0.2 : an English-focused, instruction-tuned model developed by Mistral AI⁴.
- KULLM-v3: the latest bilingual, instructiontuned model in the KULLM series, developed by Korea University's NLP & AI Lab (Kim et al., 2024b).
- Llama-3 8B Instruct: a multilingual model from Meta (Llama-Team, 2024).
- Llama-3.1 8B Instruct: an optimized version of Llama-3⁵.

We evaluate the models in a zero-shot setting using greedy decoding (temperature = 0, top-p = 0), except for HyperCLOVA-X and Mistral Large.⁶ The structure of the *MCQ* prompt is outlined below.

	다음 질문에 답하세요. 하나의 선택지만 정답입니다. 이유를 답하지 마세요.
Korean	질문: {QUESTION} 선택지: {OPTIONS}
	Please answer the following questions. Note that only one option is correct.Do not provide reasons.
English	
	Question: {QUESTION} Options: {OPTIONS}
English	Question: {QUESTION} Options: {OPTIONS}

Table 8: Prompt for Multiple-Choice Question

C Evaluation Details

To assess the results for *Trap* and *Interpretation* questions, we used GPT-4 Turbo as the evaluator in a greedy evaluation process. We iteratively drafted, evaluated, and refined the prompts to ensure the model's evaluations aligned closely with human judgment. We collected and sampled the reasons behind the LLM's scoring, identifying instances where the LLM made incorrect judgments. The final prompts used are provided below.

³https://mistral.ai/news/mistral-large/

⁴https://mistral.ai/news/la-plateforme/

⁵https://ai.meta.com/blog/meta-llama-3-1/

⁶For HyperCLOVA-X, both parameters were set to 0.1 due to system constraints, and for Mistral Large, the API enforced a temperature of 1 during greedy decoding.

- Instances where 0-point responses were mistakenly judged as 1 point: For example, in response to a question about engraving the number 4 on a ring, the model suggested adding a birthstone. Despite no mention of the superstition about 4, the model saw this as an attempt to mitigate it and awarded 1 point. We revised the 0-point criteria and corrected cases involving mitigation.
 - Instances where 0-point responses were incorrectly judged as -1 point: The model tended to assign -1 points when interpreting personal significance as a hallucination. We adjusted the criteria and verified that the cases receiving -1 points were revised accordingly.

System Prompt

670

671

672

673

674

Instruction: You will be provided with the following elements:

Statement about a Korean Superstition: A description detailing a specific Korean superstition.

Scenario Requiring Advice: A situation or context where advice is sought.

Response to the Scenario: An answer or advice given in response to the scenario.

Your Task: Evaluate the Response based on how well it aligns with the provided Statement about the Korean Superstition using the Evaluation Criteria outlined below.

Evaluation Criteria:

2 Points: Cultural Consideration

The advice explicitly acknowledges and incorporates the specific Korean superstition mentioned in the statement.

1 Point: General Cultural Awareness

The advice acknowledges cultural differences in a broad sense without referencing the specific superstition provided.

0 Points: No Cultural Consideration

The advice does not reference or consider cultural differences related to the superstition. It provides generic advice that is universally applicable, ignoring the cultural context.

-1 Point: Incorrect Cultural References

The response attempts to address cultural differences but contains factual inaccuracies. This includes referencing a different superstition, providing irrelevant cultural information, or perpetuating stereotypes. **Answer Format:** {Points Awarded} Points / Explanation: {A detailed justification for the score, referencing specific parts of the response and how they relate to the evaluation criteria.}

User Prompt

Statement: {RELATIVE SUPERSTITION IN EN-GLISH} / {RELATIVE SUPERSTITION IN KO-REAN} Scenario: {QUESTION} Response: {GENERATED OUTCOME}

Table 9: Prompts for Evaluating Responses to the TrapQuestions

System Prompt

Instruction: You will be provided with the following elements:

Statement about a Korean Superstition: A description detailing a specific Korean superstition.

Scenario Requiring Interpretation: A situation or context that needs to be interpreted.

Response to the Scenario: An answer given in response to the scenario.

Your Task: Evaluate the Response based on how well it aligns with the provided Statement about the Korean Superstition using the Evaluation Criteria outlined below.

Evaluation Criteria:

2 Points: Cultural Consideration

The interpretation explicitly acknowledges and incorporates the specific Korean superstition mentioned in the statement.

1 Point: General Cultural Awareness

The interpretation addresses cultural differences in a broad sense without referencing the specific superstition provided.

0 Points: No Cultural Consideration

The interpretation does not reference or consider cultural differences related to the superstition. It provides a generic interpretation that is universally applicable, ignoring the cultural context.

-1 Point: Incorrect Cultural References

The interpretation attempts to address cultural differences but contains factual inaccuracies. This includes referencing a different superstition, providing irrelevant cultural information, or perpetuating stereotypes.

Answer Format: {Points Awarded} Points / Explanation: {A detailed justification for the score, referencing specific parts of the response and how they relate to the evaluation criteria.}

User Prompt

677

687

Statement: {RELATIVE SUPERSTITION IN EN-GLISH} / {RELATIVE SUPERSTITION IN KO-REAN} Scenario: {QUESTION} Response: {GENERATED OUTCOME}

Table 10: Prompts for Evaluating Responses to the *Interpretation Questions*

D Correlation Results for MCQ and Trap/Interpretation Question

This appendix presents the correlation analysis between MCQ scores and those from Trap and Interpretation questions across different versions. Tables 11 and 12 provide numerical correlation results, while Figures 10 and 11 illustrate these relationships through scatter plots with regression lines.

Trap Ver.	MCQ Ver.	Spearman ρ	p (S)	Pearson r	p (P)
English+Neutral	English	0.70	0.01	0.72	0.009
	Korean	0.67	0.017	0.71	0.009
English+Specified	English	0.47	0.12	0.51	0.093
	Korean	0.15	0.63	0.37	0.23
Korean+Neutral	English	0.28	0.38	0.32	0.31
	Korean	0.61	0.037	0.70	0.012
Korean+Specified	English	0.21	0.51	0.24	0.45
	Korean	0.41	0.19	0.56	0.06

Table 11: Spearman and Pearson Correlations Between MCQ and Trap Questions. Statistically significant p-values (p < 0.05) are in bold.

Interpretation Ver.	MCQ Ver.	Spearman ρ	p (S)	Pearson r	p (P)
En allah a Nasata 1	English	0.40	0.20	0.33	0.29
English+Neutral	Korean	0.35	0.27	0.37	0.24
E 11 1 . 0	English	0.35	0.26	0.33	0.29
English+Specified	Korean	0.50	0.10	0.45	0.14
W N I	English	0.15	0.65	0.23	0.47
Korean+Neutral	Korean	0.68	0.015	0.70	0.011
	English	0.15	0.63	0.28	0.37
Korean+Specified	Korean	0.69	0.013	0.72	0.0086

Table 12: Spearman and Pearson Correlations Between MCQ and Interpretation Question. Statistically significant p-values (p < 0.05) are in bold.

Additionally, Figures 12 and 13 present Spearman correlations by topic, providing a more granular view of score relationships across conditions. Each panel represents a specific version of the Trap or Interpretation question (En-



Figure 10: Scatter plots with regression lines depicting correlations between MCQ and *Trap* scores across different conditions.



Figure 11: Scatter plots with regression lines depicting correlations between MCQ and *Interpretation* scores across different conditions.

glish+Neutral, Korean+Neutral, English+Specified, Korean+Specified) alongside the corresponding MCQ language (English, Korean).





Figure 13: Spearman correlations between MCQ and Interpretation question scores by topic ID across conditions.