Can Code-Switched Texts Activate a *Knowledge Switch* () in LLMs? A Case Study on English-Korean Code-Switching

Anonymous ACL submission

Abstract

001 Recent large language models (LLMs) demonstrate multilingual abilities, yet they are English-002 003 centric due to dominance of English in training corpora. The limited resource for low-resource languages remains a crucial challenge. Codeswitching (CS), a phenomenon where multilin-006 gual speakers alternate between languages in 007 800 a discourse, can convey subtle cultural and linguistic nuances that can be otherwise lost in 009 translation and elicits language-specific knowl-010 edge in human communications. In light of this, 011 we investigate whether code-switching can ac-012 013 tivate, or identify and leverage knowledge for reasoning when LLMs solve low-resource lan-014 guage tasks. To facilitate the research, we first 015 present ENKOQA, a synthetic English-Korean 016 CS question-answering dataset. We provide 017 comprehensive analysis on a variety of multilin-018 gual LLMs by subdividing activation process 019 into knowledge identification and knowledge leveraging. Our results demonstrate that com-021 pared to English text, CS can faithfully activate 022 knowledge inside LLMs especially on language-024 specific domains, suggesting the potential of code-switching on low-resource language tasks. 025

1 Introduction

027

029

030

031

037

039

040

041

Large language models (LLMs) have continuously evolved through time to exhibit advanced multilingual capabilities, enabled by training on massive datasets that include text in many different languages. However, these sources are typically skewed toward English, creating an inconsistent performance across different languages (Chen et al., 2024; Zhang et al., 2024). The limited availability for real-world user queries in low-resource languages remains a crucial challenge for achieving robust multilingual models. Prior works attempt to mitigate this issue through machine translation (Artetxe et al., 2023; Bareiß et al., 2024), but crucial semantic nuances may be lost in translation, and machine translation errors are inevitable.



Figure 1: A motivating example of knowledge identification between languages. Compared to a question in English (*top*), a bilingual speaker can "activate" more relevant knowledge with a question in CS (*bottom*).

042

043

044

045

046

047

048

049

050

051

052

053

054

055

056

057

058

059

In human multilingual societies, code-switching (CS), or the practice of alternating between two or more languages within an utterance, is used to fill in lack of language proficiency, to emphasize certain emotions or points, or for group identity (Heredia and Altarriba, 2001). Moreover, code-switching functions as an effective tool to embed cultural meanings. Expressing certain concepts in original language can convey subtle cultural and linguistic nuances that can be lost in translation, and knowledge related to certain language are more likely to be more memorized in its own language. As shown in Figure 1, when a human English-Korean bilingual is given a question that is closely related to Korean culture, a question in English and Korean code-switching is more capable of recalling knowledge about "몽유도원도"¹, because the concept is more familiar in Korean than in English.

¹A landscape painting by An Gyeon in the early Joseon Dynasty requested by Prince Anpyeong, after his dream about Shangri-la. The painting is drawn on silk with ink.

This observation raises intriguing insight about the impact of code-switching in multilingual societies and the potential for equivalent effect in LLMs. Given that code-switching facilitates target language-specific knowledge in human communications, we investigate whether the same applies to English-centric LLMs when solving low-resource language tasks. Therefore, we ask ourselves the following research question: **Can code-switched texts activate language-specific knowledge, or turn on a "knowledge switch" in LLMs?** By *knowledge activation*, we refer to the overall process of identifying what knowledge is required, and applying knowledge to answer the question.

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

078

079

081

082

083

084

085

087

089

090

091

093 094

095

096

097

099

100

101

102

103

104

105

106 107

108

109

110

111

To answer the question, we subdivide knowledge activation process into two tasks: (1) In *Knowledge Identification* task, we investigate if querying LLMs in CS and English yield different knowledge from its encoded memory. Specifically, we evaluate the quality of knowledge from different linguistic settings in terms of faithfulness and helpfulness. (2) In *Knowledge Leveraging* task, we observe if LLMs can faithfully ground on identified knowledge for solving question-answering (QA) task.

There have been continuous, if not abundant, researches on code-switching in the field of computational linguistics (Aguilar et al., 2020; Rizvi et al., 2021). Recently, after the emergence of LLMs with impressive multilingual abilities, a line of work have discovered LLMs' abilities in CS (Huzaifah et al., 2024; Yong et al., 2023; Zhang et al., 2023a). However, the focus of such works are only limited to understanding and generating CS of LLMs, while the effectiveness of CS in tasks that involve lowresource language has not yet been explored. To the best of our knowledge, this work is the first to comprehensively analyze the effectiveness of codeswitching on knowledge activation to LLMs.

Meanwhile, a crucial challenge when it comes to code-switching is the data scarcity. There is a limited number of CS datasets, let alone culturefocused data (Doğruöz et al., 2021). Since CS often happens in conversations, data are not easily available and the quality is not ensured. To address the shortage of data, efforts have been made to synthetically generate code-switching corpus based on linguistic theories (Pratapa et al., 2018; Rizvi et al., 2021; Salaam et al., 2022). However, these works rely on syntactic parsers and part-of-speech taggers that support limited languages, and the quality of text are highly dependent on the performances of those tools. Therefore, we first construct ENKOQA, a synthetic English-Korean code-switching dataset to explore the potential of CS in low-resource language task. Following Matrix Language Frame Model (Myers-Scotton, 1997), we synthesize Korean QA datasets (Kim et al., 2024b; Son et al., 2024) that encompass various aspects of Korea into English-Korean code-switched questions. 112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

We conduct experiments with ENKOQA and provide extensive analysis on a wide range of multilingual LLMs. The experimental results reveal that CS is able to faithfully activate language-specific knowledge that are encoded in multilingual LLMs compared to high-resource language and target language translation; this tendency was more prominent on domains that specifically requires knowledge in target language and culture.

The contributions of our work are as follows:

- To the best of our knowledge, this work is the first to comprehensively analyze the effectiveness of code-switching on knowledge activation to LLMs by introducing two tasks.
- We propose a qualified English-Korean codeswitching QA dataset that is synthesized upon two Korean-centric datasets, and conduct extensive experiments on various families of multilingual LLMs.
- Experimental results on extensive LLMs indicate that code-switching has advantages in knowledge activation especially on languagespecific domains, suggesting the potential of code-switching text as a tool for conveying cultural nuances in target language tasks.

2 Preliminaries & Related Work

In this section, we provide preliminary knowledge about code-switching, and explore relevant studies from conventional and computational linguistics.

2.1 Code-Switching Theories

Many linguistic theories attempt to explain the grammatical construction of code-switched text, such as Equivalence Constraint (EC) theory and Free Morpheme Constraint (FMC) theory proposed by Poplack (1980). EC theory suggests that code-switching occurs at points in a sentence where the structures of both languages are grammatically compatible. FMC theory suggests that code-switching cannot occur between a bound morpheme and a lexical base. (*e.g.*, "He is look-ando for a book." is a wrong code-switch.)

However, these theories have limitations in that the theory can only be applied to two language with similar or equivalent syntactic structures. EC and FMC theories are not applicable to English-Korean code-switching text, due to the different sentence structure of Korean and English (Park and Yun, 2021). In this regard, we adopt Matrix Language Frame Model to construct our codeswitching dataset.

Matrix Language Frame Model 2.2 169

160

161

162

163

164

165

166

167

168

187

191

197

201

205

206

208

Matrix Language Frame (MLF) model is a code-170 switching theory proposed by Myers-Scotton 171 (1997). MLF model posits that in any instance of 172 code-switching, one language provides the morpho-173 syntactic framework of the sentence. This is known 174 as the matrix language. The other language, called 175 the embedded language, contributes to additional 176 content, usually in the form of words or phrases, 177 but follows the grammatical rules set by the matrix 178 language. In other words, matrix language domi-179 nates the sentence structure, while the embedded 180 language is integrated within that structure. Content 181 morphemes can be in both languages, but functional 182 morphemes come from matrix language. Taking 183 Figure 1 as an example, "그 내용" which translates 184 to "its contents" can be embedded into English sen-185 tence, but functional morpheme such as "to" cannot. 186

2.3 Code-Switching for Language Models

Previous works introduce benchmarks for evaluat-188 ing code-switching ability of multilingual language 189 models across multiple tasks (Aguilar et al., 2020; 190 Khanuja et al., 2020). More recent works focus on 192 the capability of LLMs in code-switching. Zhang et al. (2023a) discover performance of multilin-193 gual LLMs in various code-switching tasks, includ-194 ing sentiment analysis and language identification. 195 Yong et al. (2023) explore prompting multilingual 196 LLMs to generate code-mixed data. Shankar et al. (2024) introduce a prompting technique called in-198 context mixing for effective in-context learning in 199 LLMs. Although these benchmarks encompass a va-200 riety of tasks, the analysis of LLMs' code-switching capabilities in terms of knowledge retrieval and uti-202 lization has not yet been investigated. 203

2.4 Code-Switched Data Synthesis

Data synthesis for code-switching has been approached in various ways. Several studies utilize parsers and neural models to synthesize codeswitched text based on EC theory (Pratapa et al.,

2018; Rizvi et al., 2021). Similarly, Salaam et al. 209 (2022) extract phrases from source language and 210 reintegrate them into target language. In recent ef-211 forts to address data scarcity in low-resource set-212 tings, LLMs have been employed to generate syn-213 thetic data (Li et al., 2023). However, using LLMs 214 specifically for synthesizing code-switched data re-215 mains unexplored. 216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

ENKOQA: English-Korean 3 **Code-Switching QA Testset**

To compare the effectiveness of code-switching with dominant language and translation in target language when performing language-specific tasks, we introduce ENKOQA, a synthetic English-Korean code-switching dataset that is designed based on MLF model. In this section, we first discuss the details of data construction (\S 3.1), and evaluate performances of LLMs on the dataset (\S 3.2, 3.3).

3.1 Dataset Construction

Data Sources. We leverage two multiple-choice Korean-centric question-answering datasets that encompass various aspects of Korean language and culture. CLIcK (Kim et al., 2024b) is a Korean benchmark dataset designed to test Korean cultural and linguistic knowledge collected from various official Korean exams and textbooks, e.g., College Scholastic Ability Test of Korea (CSAT). HAE-RAE (Son et al., 2024) is a Korean benchmark dataset originally crafted to capture cultural and contextual nuances inherent to the Korean language, sourced from official Korean exams, textbooks, and text on the internet. In this work, we focus on categories about Korean society to evaluate the effect of CS on activating Korean-specific knowledge. Specifically, we collect 1,995 pairs of eight categories from CLIcK, and 1,027 pairs of five categories from HAE-RAE, resulting in 2,372 QA pairs in nine categories: Popular, Economy, Politics, Tradition, General Knowledge, Society, Geography, History, and Law. More details of original datasets are provided in Appendix A.1.

Automatic Translation. As most LLMs are trained on English-dominant corpora, we regard the English-centric LLM as a bilingual whose matrix language is English but also fairly competent in Korean. To generate code-switched text that follows the MLF model, we need parallel data in Korean and English to extract semantically important words or phrases from Korean text and embed

Model		Economy	General	Geography	History	Law	Politics	Popular	Society	Tradition	Total
	CS	91.53	78.41	69.04	74.79	55.86	90.48	95.12	63.70	85.14	78.23
GPT-40	EN	89.83	75.00	66.19	61.97	52.64	84.52	95.12	60.40	74.32	73.33
	KO_t	89.83	71.59	60.14	63.03	48.74	85.71	92.68	56.44	75.23	71.49
	CS	71.19	47.73	44.48	32.91	35.40	70.24	80.49	49.17	57.21	54.31
GPT-3.5	EN	71.19	48.86	45.55	36.32	36.55	66.67	63.41	52.64	62.61	53.76
	KO_t	62.71	26.70	31.67	26.71	24.83	48.81	58.54	37.79	49.55	40.81
	CS	93.22	72.16	72.95	73.08	62.53	86.90	95.12	67.66	84.23	78.65
Claude 3.5	EN	89.83	71.59	67.97	61.54	55.63	85.71	92.68	63.20	75.23	73.71
	KO_t	64.41	47.73	54.09	54.49	45.52	69.05	82.93	52.31	61.71	59.14
	CS	83.05	55.11	54.09	63.46	42.76	80.95	85.37	54.29	75.23	66.03
Solar	EN	74.58	46.02	49.47	39.53	42.76	77.38	65.85	51.16	62.61	56.60
	KO_t	81.36	50.57	56.94	58.12	46.44	82.14	78.05	54.95	70.27	64.31
	CS	79.66	51.70	50.53	49.36	44.14	80.95	75.61	57.43	65.77	61.68
Llama3 70B	EN	83.05	57.39	50.53	45.94	45.75	73.81	73.17	53.30	66.67	61.07
	KO_t	76.27	50.57	46.98	43.80	38.16	70.24	82.93	51.49	61.26	57.97
	CS	69.49	40.34	36.30	35.68	35.63	75.00	73.17	45.05	54.05	51.63
Llama3 8B	EN	64.41	39.77	37.72	37.39	32.64	67.86	63.41	45.21	53.60	49.11
	KO_t	61.02	38.07	38.79	32.48	33.33	65.48	65.85	43.73	50.90	47.74
	CS	79.66	46.02	48.75	41.03	45.29	77.38	78.05	54.79	65.32	59.59
Gemma2 27B	EN	84.75	53.41	48.40	40.60	41.84	72.62	78.05	54.95	63.96	59.84
	KO_t	77.97	44.89	44.84	41.67	41.84	73.81	75.61	50.66	59.46	56.75
	CS	79.66	42.05	44.13	40.17	41.15	73.81	80.49	53.30	65.77	57.84
Gemma2 9B	EN	76.27	46.02	49.47	38.46	42.30	69.05	73.17	52.15	63.51	56.71
	KO _t	76.27	42.05	41.99	34.62	40.23	71.43	82.93	51.98	58.11	55.51

Table 1: QA performances of multilingual LLMs on CS, English, and translated Korean settings. **Bold** indicates the highest score among the three baselines from each model. **Green** indicates the highest score from each domain.

into English text. We first automatically translate all Korean query-choices pairs into English using gpt-3.5-turbo, where the model is instructed to translate the query and choices to English with an one-shot demonstration. Lastly, human supervision was done to ensure translation quality.

258

259

260

261

262

263

Generating Candidates in Different Levels. 264 Now that we obtain parallel data in both languages, 265 the next step is to embed Korean content mor-266 phemes into English sentence. As code-switching 267 268 mostly happens spontaneously, there does not exist a certain formula for mixing two languages. More-269 over, replacing every content word with its Korean 270 equivalent may seem rather artificial. To address 271 this, we simulate a natural code-switching by creat-272 ing various versions of code-switched texts at dif-273 ferent ratios (30, 50, 70, and 90%), then selecting a 274 version that represents the best quality and most nat-275 uralness. Specifically, given a question in both lan-276 guages and a specified proportion, gpt-3.5-turbo 277 identifies content words from the Korean question 278 and integrates them into the English question ac-279 cording to the specified proportion. To collect con-281 texts of various semantic importance, we employ two prompts that define "content word" differently; 282 one defines content words as noun phrases, while 283 the other identifies them as semantically important 284 elements within the context. Eight code-switched 285

candidates are collected per question, from which human annotators select a single candidate that most faithfully follows MLF structure. Comprehensive details about dataset construction are provided in Appendix A. 286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

3.2 Experimental Settings

Models. We conduct extensive analysis on two groups of state-of-the-art multilingual LLMs: (1) Proprietary LLMs that are available via APIs, such as GPT-3.5, GPT-40 (OpenAI, 2023), and Claude 3.5 Sonnet (Anthropic, 2024). (2) Open-source LLMs such as Solar (10.7B, Kim et al., 2024a), Llama3 (8B, 70B, Dubey et al., 2024), and Gemma2 (9B, 27B, Gemma Team, 2024). More details about the models are in Appendix B.1.

Baselines. To compare performances of LLMs in various language settings, we evaluate on CS, English, and translated Korean (KO_i) questions. Korean translation baseline simulates cases where machine translation is adopted to convert task data from English into the target language in lowresource language tasks. We back-translate English translation text to Korean using gpt-3.5-turbo. Prompts that we used for inference are provided in Table 13. We also conduct experiments on the original Korean questions, but do not consider it as major baseline, because we aim to examine the effect of

403

404

405

358

359

code-switching compared to dominant language,
rather than demonstrate the performance of lowresource language. Please refer to Appendix B.3 for
further discussions.

317 **3.3 Results**

318Overall. As shown in Total column from Table 1,319the performance on CS significantly outperforms320English and KO, across most LLMs in average. The321gap between CS and other baselines is especially322prominent in GPT-40 and Claude 3.5, where CS323peaks in all domains.

CS questions excel at language-specific domains. 324 While CS outperforms other baselines in many 325 domains, it is worth noting that the gap between 326 CS and English is substantially large on language-327 sensitive domains such as History and Tradition, 328 both of which target language is essential for pre-329 serving information or terminology. Even Llama3 and Gemma2 models which relatively do not per-331 form well on CS questions, show higher scores on 332 CS for such domains. On the other hand, the phe-333 nomenon is less consistent for general knowledge 334 (e.g., Society, General), and domains that require 335 expert-level knowledge (e.g., Politics, Law). 336

CS surpasses translated Korean on most models. 337 We compare code-switching with translated Korean 338 339 translation to observe whether CS has advantages in minimizing translation errors. Except for Solar, 340 KO_t generally shows lowest performance among 341 three baselines. This suggests that while translating 342 task in target language is not the best practice, CS 343 can faithfully encapsulate meanings and linguistic 344 cues that may be lost in translation, highlighting 345 the potential of leveraging CS for performing non-346 dominant language tasks. 347

Ratios do not affect performance. To ensure that the ratio of code-switching does not influ-349 ence models' performances and our dataset is con-350 structed under fair process, we calculate Code-351 Mixing Index (CMI) scores (Srivastava and Singh, 352 353 2021) and report corresponding accuracy in Tradition and History domains. As shown in Table 3, 354 we can see that accuracy scores are quite evenly 355 distributed across all ratios, suggesting that there is no distinct tendency between CMI and accuracy. 357

4 Can Code-Switched Questions Activate a "Knowledge Switch" in LLMs?

From Section 3.3, we observe that most LLMs are able to answer correctly to questions in CS than in other baselines. To further investigate on the effectiveness of CS in activating language-specific knowledge, we formulate two tasks: *Knowledge Identification* and *Knowledge Leveraging*. We evaluate the tasks in CS and English questions, the two baselines that share the same matrix language.

4.1 Knowledge Identification

Task Description. When a human English-Korean bilingual is given a question about Korean culture, they will first try to identify what specific knowledge is required to answer the question, and then apply the knowledge to find the correct answer. Depending on which language the question is written in, the quantity and quality of the knowledge may vary, as described in Figure 1. Languagespecific knowledge is likely to be encoded much abundantly in its own language, so reading the question in CS will allow more effective knowledge activation than in English. In this sense, knowledge identification task evaluates LLMs' ability to identify what knowledge is prerequisite for the question. Specifically, the LLM is asked to write a list of factual knowledge that are necessary for solving the given question in one or two sentences.

Evaluation Criteria. For a qualitative analysis on knowledge identification, we evaluate the quality of a knowledge list based on two criteria: *Faithfulness* evaluates whether the generated knowledge is factually correct and the model does not output hallucination. *Helpfulness* evaluates whether the knowledge is relevant to the question, and helpful for answering the question correctly.

4.2 Knowledge Leveraging

Task Description. We refer to Knowledge Leveraging as applying the identified knowledge into reasoning. In specific, the model should be able to find a correct answer based on the knowledge it has identified from the Knowledge Identification task. Therefore, we provide knowledge identified by each model and instruct the model to find the answer using the knowledge. To encourage the models to properly ground on knowledge, we adopt Chainof-Thought reasoning (Wei et al., 2023) and prompt the models to generate reasoning steps that lead to



Figure 2: Human evaluation results on faithfulness (*top*) and helpfulness (*bottom*) of knowledge lists identified from **CS** questions and **English** questions.

the final answer. We conduct experiments on the entire dataset and report accuracy score.

4.3 Experimental Setup

406

407

408

Implementation Details. We conduct experi-409 ments with the same models as in Section 3.2. For 410 knowledge identification, we instruct the model to 411 write a list of factual knowledge that are required for 412 solving the given question in one or two sentences. 413 For knowledge leveraging, we pass on previously 414 identified knowledge and ask the model to select an 415 answer and explain why. The full-length prompts 416 are provided in Table 14 and 15. 417

Evaluating Knowledge Identification. In order 418 to effectively evaluate knowledge identification re-419 420 sults, we refer to Section 3.3 and choose two domains where CS performance is higher (*i.e.*, History, 421 Tradition), and two domains that have minimum 422 difference (*i.e.*, General, Law). Moreover, we select 423 four models with different performances and sizes 424 (i.e., GPT-40, Solar, Gemma2 27B, Gemma2 9B). 425 Specifically, we sample 10 questions from each do-426 main and model, resulting in 160 samples. Then, 427 we conduct human and LLM-based evaluation on 428 identified knowledge. 429

Human Evaluation We employ four human evaluators who are fluent in both Korean and English
and completed Korean public education, thus qualified to evaluate questions sourced from Korean
proficiency tests for foreigners and the Korean Col-



Figure 3: Human evaluation results on pairwise comparisons between knowledge lists identified from **CS** questions and **English** questions.

lege Scholastic Ability Test. For faithfulness and helpfulness, the evaluator is asked to rate a knowledge list on a Likert scale from 1 to 3. In pairwise evaluation, we provide two knowledge lists in a random order and ask the evaluator to select a list that is overall more effective for answering the question. Details on evaluation criteria and evaluator information are provided in Appendix C.1 and C.2. 435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

LLM-based Evaluation As we conduct human evaluation on quite small amount of samples, we additionally conduct LLM-as-a-judge evaluation (Zheng et al., 2023) to amplify our analysis. Specifically, we use GPT-40 as the evaluator, using identical instructions with human evaluators on 40 questions for 9 domains and 8 models, 360 samples in total. Full prompts are provided in Appendix C.1.

5 Analysis on Knowledge Identification

5.1 Human Evaluation

Faithfulness. In the upper row of Figure 2, we 453 observe a significant gap in faithfulness scores be-454 tween CS and English in both History and Tradition. 455 The discrepancy is more salient in Tradition where 456 cultural nuances is much important, implying that 457 asking questions in CS is much successful in cap-458 turing cultural nuances and meanings. In General 459 domain, the scores for CS and English are almost 460 identical (or even better in English for Gemma2 461 9B), indicating that the difference in knowledge 462 activated by CS questions compared to English 463 questions is minimal when addressing general and 464 common facts. In Law, although knowledge from 465 CS is slightly more faithful than that from English, 466 their absolute scores are lower than those in other 467 domains, suggesting that models fail to identify 468 faithful knowledge that requires domain expertise. 469



Figure 4: LLM-as-a-judge evaluation results on pairwise comparison between knowledge lists identified from **CS** questions and **English** questions.

Helpfulness. The lower row of Figure 2 presents 470 evaluation results for helpfulness. It is intuitive that 471 faithful knowledge serves as a valuable source for 472 answering questions, and as a result, the evalua-473 tion of helpfulness shows a similar trend to that 474 of faithfulness. In History and Tradition, the gap 475 between CS and English becomes larger in helpful-476 ness, emphasizing the effectiveness of the CS set-477 ting in identifying both faithful and helpful knowl-478 479 edge. It is also notable that the scores for helpfulness are particularly high for GPT-40 and Solar, 480 models in which performance in CS surpasses that 481 in English to a large extent (\S 3.3). In contrast, the 482 helpfulness scores in the Law domain are consid-483 erably lower for both CS and English compared to 484 other domains. Given that the Law domain requires 485 expert-level legal knowledge, the models struggle 486 to grasp the legal context, leading to difficulties in 487 accurately identifying helpful knowledge sources 488 from both CS and English questions. 489

Pairwise Comparison. In Figure 3, the win ratio 490 for CS is higher in History and Tradition, demon-491 strating that CS questions can activate more essen-492 tial knowledge sources for question answering. On 493 the contrary, in domains where CS does not show 494 its effectiveness, the win ratio of CS is compara-495 496 tively lower (*i.e.*, General) or the ratio of Tie is high (*i.e.*, Law). Especially in the case of Law, the qual-497 ity of knowledge lists generated from CS questions 498 is evaluated as equivalent to, or even worse than, 499 that generated from English questions. 500

5.2 LLM-based Evaluation

We observe in Figure 6 and Figure 7 that the score gap between CS and English in both faithfulness and helpfulness are minimal. In fact, CS scores are even or lower for some cases, which are inconsistent with human evaluation results. However, it is still worth noting that LLM-as-a-judge also assigns higher scores for advanced models, and overall scores were lower in History and Tradition. 501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

On the other hand, LLM judgement scores in pairwise evaluation generally agree with the human evaluations. We compute Cohen's Kappa (κ) score in Table 6, and follow interpretations from Landis and Koch (1977).² Consistent with human evaluation, the LLM judge votes CS for most cases, and the agreement is stronger with advanced models (*i.e.*, GPT-40), on culture-intensive domains (*i.e.*, History, Tradition).

While other domains fairly agree with human judgment, Law shows exceptional results. Specifically, the LLM-as-judge evaluation reports a significantly higher win ratio for CS in the Law domain compared to human evaluation. However, considering that tie ratio is substantial in human evaluation as well, we speculate that LLM-as-a-judge gives a win to CS on knowledge that human evaluators regarded comparable quality with English setting.

6 Analysis on Knowledge Leveraging

We present the visualized results of accuracy in both CS and English settings in Figure 5, with detailed scores reported in Table 7.

Main Observations. Consistent with the results in Section 3.3, all models demonstrate generally higher performances for CS questions compared to English questions. The results indicate that CS effectively activates knowledge across various domains while activating in dominant English language is suboptimal. GPT-40, Claude, and Solar exhibit higher CS performance than English across all domains. These models not only identify faithful and helpful knowledge (§ 5.1), but also answer questions by accurately grounding on that knowledge; this shows that CS questions robustly activate essential knowledge in these models. On the contrary, Llama3 and Gemma2 families show poor performance in both CS and English settings in several domains, such as Geography and Law. Taking into

 $^{^{2}}$ Landis and Koch (1977) interprets 0–0.20 as slight, 0.21–0.40 as fair, 0.41–0.60 as moderate, 0.61–0.80 as substantial, and 0.81–1 as almost perfect agreement.



Figure 5: Radar charts of knowledge leveraging performances on all domains across various multilingual LLMs. **Green** line is code-switching and dashed **gray** line is English. We report accuracy for the evaluation metric.

account that these domains require domain-specific expertise, it is likely that their lack of understanding contributes to low accuracy, let alone CS failing to activate Korea-focused knowledge.

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

Knowledge Identification and Leveraging both matters. We demonstrate that qualified knowledge identification is prerequisite for knowledge activation of CS. The win ratio of History knowledge by GPT-3.5 was relatively poor compared to others (Figure 4), leading GPT-3.5 to be the only model that did not benefit from CS. Similarly in Law, Figure 2 and 3 show that helpfulness and pairwise scores for knowledge by Gemma2 9B and 27B are lower than others, which are responsible for their suboptimal performance in CS. We provide further qualitative case study in Appendix D.4.

English questions hallucinate more than CS. 564 Although we informed the models that the answer 565 is in one of the choices, we notice that the ma-566 jority of incorrect responses were "None of the 567 above". The errors may derive from either hallucinated knowledge or failing to follow instructions 569 faithfully. Therefore, we provide additional analy-570 sis on erroneous outputs in Table 8. We report the 571 results in the format of number of errors that de-572 573 rived from knowledge hallucination / total number of None errors. Errors that are not from halluci-574 nation are caused by poor instruction-following. 575 Overall, we observe that performance on English 576 questions results in more errors compared to CS 577

across all LLMs, and most of them were hallucination errors. This indicates that models hallucinate much frequently when English questions are given, again highlighting the effectiveness of CS over English setting. It is also worth noting that Gemma2 families hallucinate largely on History and General, supporting our finding in Figure 2 and Figure 5 which respectively illustrates poor performance on human evaluation and QA accuracy. 578

579

580

581

582

583

584

585

586

587

7 Conclusion

We explore the efficacy of code-switching in acti-588 vating language-specific knowledge embedded in 589 LLMs. Utilizing two Korean-centric QA datasets, 590 we synthesize ENKOQA, a qualified English-591 Korean code-switching QA dataset and conduct 592 experiments on various multilingual LLMs. Our 593 analyses demonstrate that LLMs can simulate a 594 similar code-switching effect with human commu-595 nications of facilitating low-resource knowledge 596 within LLMs, particularly in language-specific do-597 mains. Regarding this finding, we suggest that code-598 switching can be an effective strategy for solving 599 low-resource language tasks. Also, augmenting low-600 resource datasets into code-switching text can am-601 plify resource and mitigate data scarcity challenge. 602 We hope our work can motivate NLP community 603 to explore more potential of code-switching for de-604 veloping robust multilingual LLMs. 605

606 Limitations and Future Work

In this work, we focus on code-switching between 607 English and Korean, specifically limiting the scope 608 to Korea-specific knowledge. However, it is impor-609 tant to note that this study serves as a single case 610 focused on the Korean context and leaves room for 611 expanding the scope of code-switching to other cul-612 tures and languages. For future research, we aim to 613 investigate whether the knowledge activation effect 614 also occurs in other language settings. 615

Another limitation of our work is that we con-616 duct human evaluations on only a subset of LLMs, 617 domains, and questions. Evaluating the quality (i.e., 618 faithfulness and helpfulness) of knowledge in code-619 switched text presents inherent and practical chal-620 lenges, as it necessitates evaluators to be fluent bilin-621 guals. Consequently, we present only partial results 622 for the knowledge identification task. 623

Lastly, as we rely on a LLM, specifically 624 gpt-3.5-turbo, to synthesize our code-switching 625 dataset, the performance of the LLM can affect the 626 quality of the dataset. To mitigate the risk of er-627 roneous samples and to fully leverage the LLM's 628 capabilities, we engage reliable human annotators 629 to review the samples and verify their quality. Also, 630 631 as we formulate our code-switching dataset with gold English and Korean, in a more realistic sce-632 nario where a monolingual English speaker creates 633 code-switching text, sentences would be created 634 automatically without any additional supervision. 635

> In the future, we aim to investigate more potential of code-switching in diverse aspects, including instruction-tuning of LLMs to users effectively using code-switching for multilingual tasks. As we have demonstrated synthesizing monolingual datasets into code-switching text, we hope our work can inspire NLP community to explore the capability of code-switching in enhancing and utilizing multilingual LLMs.

645 Ethical Consideration

636

637

638

639

640

641

642

643

644

646Our work utilizes large language models for data647construction. Recent work has highlighted the risks648of LLMs in hallucination (Zhang et al., 2023b). In649order to prevent any hallucination or harmful con-650tents, we ensure that human annotators examined651each sample carefully and create dataset safely.

References

Gustavo Aguilar, Sudipta Kar, and Thamar Solorio. 2020. LinCE: A centralized benchmark for linguistic code-switching evaluation. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 1803–1813, Marseille, France. European Language Resources Association. 652

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

- Anthropic. 2024. Claude 3.5 sonnet.
- Mikel Artetxe, Vedanuj Goswami, Shruti Bhosale, Angela Fan, and Luke Zettlemoyer. 2023. Revisiting machine translation for cross-lingual classification. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6489–6499, Singapore. Association for Computational Linguistics.
- Patrick Bareiß, Roman Klinger, and Jeremy Barnes. 2024. English prompts are better for nli-based zero-shot emotion classification than target-language prompts. In *Companion Proceedings of the ACM Web Conference 2024*, WWW '24, page 1318–1326, New York, NY, USA. Association for Computing Machinery.
- Pinzhen Chen, Shaoxiong Ji, Nikolay Bogoychev, Andrey Kutuzov, Barry Haddow, and Kenneth Heafield. 2024. Monolingual or multilingual instruction tuning: Which makes a better alpaca. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 1347–1356, St. Julian's, Malta. Association for Computational Linguistics.
- A. Seza Doğruöz, Sunayana Sitaram, Barbara E. Bullock, and Almeida Jacqueline Toribio. 2021. A survey of code-switching: Linguistic and social perspectives for language technologies. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 1654–1666, Online. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The Ilama 3 herd of models. *Preprint*, arXiv:2407.21783.
- Gemma Team. 2024. Gemma 2: Improving open language models at a practical size.
- Roberto Heredia and Jeanette Altarriba. 2001. Bilingual language mixing: Why do bilinguals code-switch? *Current Directions in Psychological Science - CURR DIRECTIONS PSYCHOL SCI*, 10:164–168.
- Muhammad Huzaifah, Weihua Zheng, Nattapol Chanpaisit, and Kui Wu. 2024. Evaluating code-switching translation with large language models. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 6381– 6394, Torino, Italia. ELRA and ICCL.

814

815

816

817

818

819

820

763

Simran Khanuja, Sandipan Dandapat, Anirudh Srinivasan, Sunayana Sitaram, and Monojit Choudhury. 2020. GLUECoS: An evaluation benchmark for codeswitched NLP. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 3575–3585, Online. Association for Computational Linguistics.

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

- Dahyun Kim, Chanjun Park, Sanghoon Kim, Wonsung Lee, Wonho Song, Yunsu Kim, Hyeonwoo Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, Changbae Ahn, Seonghoon Yang, Sukyung Lee, Hyunbyung Park, Gyoungjin Gim, Mikyoung Cha, Hwalsuk Lee, and Sunghun Kim. 2024a. Solar 10.7b: Scaling large language models with simple yet effective depth upscaling. *Preprint*, arXiv:2312.15166.
- Eunsu Kim, Juyoung Suk, Philhoon Oh, Haneul Yoo, James Thorne, and Alice Oh. 2024b. CLIcK: A benchmark dataset of cultural and linguistic intelligence in Korean. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 3335–3346, Torino, Italia. ELRA and ICCL.
- J. Richard Landis and Gary G. Koch. 1977. The measurement of observer agreement for categorical data. *Biometrics*, 33(1):159–174.
- Zhuoyan Li, Hangxiao Zhu, Zhuoran Lu, and Ming Yin. 2023. Synthetic data generation with large language models for text classification: Potential and limitations. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 10443–10461, Singapore. Association for Computational Linguistics.
- Felipe Maia Polo, Lucas Weber, Leshem Choshen, Yuekai Sun, Gongjun Xu, and Mikhail Yurochkin. 2024. tinyBenchmarks: evaluating LLMs with fewer examples. In Proceedings of the 41st International Conference on Machine Learning, volume 235 of Proceedings of Machine Learning Research, pages 34303–34326. PMLR.
- C. Myers-Scotton. 1997. *Duelling Languages: Grammatical Structure in Codeswitching*. Clarendon Press.
- OpenAI. 2023. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Lorenzo Pacchiardi, Lucy G. Cheke, and José Hernández-Orallo. 2024. 100 instances is all you need: predicting the success of a new llm on unseen data by testing on a few instances. *Preprint*, arXiv:2409.03563.
- Chanjun Park, Hyeonwoo Kim, Dahyun Kim, Seonghwan Cho, Sanghoon Kim, Sukyung Lee, Yungi Kim, and Hwalsuk Lee. 2024. Open ko-llm leaderboard: Evaluating large language models in korean with koh5 benchmark. In *The 62nd Annual Meeting of the As*sociation for Computational Linguistics (ACL 2024).

- Eunsun Park and Hongoak Yun. 2021. The grammatical constraint and grammatical encoding of Korean-English code switching. *The Journal of Mirae English Language and Literature*, 26(1):177–204.
- Shana Poplack. 1980. Sometimes i'll start a sentence in spanish y termino en español: toward a typology of code-switching 1. *Linguistics*, 18:581–618.
- Adithya Pratapa, Gayatri Bhat, Monojit Choudhury, Sunayana Sitaram, Sandipan Dandapat, and Kalika Bali. 2018. Language modeling for code-mixing: The role of linguistic theory based synthetic data. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1543–1553, Melbourne, Australia. Association for Computational Linguistics.
- Mohd Sanad Zaki Rizvi, Anirudh Srinivasan, Tanuja Ganu, Monojit Choudhury, and Sunayana Sitaram. 2021. GCM: A toolkit for generating synthetic codemixed text. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 205–211, Online. Association for Computational Linguistics.
- Cesa Salaam, Franck Dernoncourt, Trung Bui, Danda Rawat, and Seunghyun Yoon. 2022. Offensive content detection via synthetic code-switched text. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6617–6624, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Bhavani Shankar, Preethi Jyothi, and Pushpak Bhattacharyya. 2024. In-context mixing (ICM): Codemixed prompts for multilingual LLMs. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4162–4176, Bangkok, Thailand. Association for Computational Linguistics.
- Guijin Son, Hanwool Lee, Suwan Kim, Huiseo Kim, Jae cheol Lee, Je Won Yeom, Jihyu Jung, Jung woo Kim, and Songseong Kim. 2024. HAE-RAE bench: Evaluation of Korean knowledge in language models. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 7993–8007, Torino, Italia. ELRA and ICCL.
- Vivek Srivastava and Mayank Singh. 2021. Challenges and limitations with the metrics measuring the complexity of code-mixed text. In *Proceedings of the Fifth Workshop on Computational Approaches to Linguistic Code-Switching*, pages 6–14, Online. Association for Computational Linguistics.
- Rajan Vivek, Kawin Ethayarajh, Diyi Yang, and Douwe Kiela. 2024. Anchor points: Benchmarking models with much fewer examples. In *Proceedings of the* 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1576–1601, St. Julian's, Malta. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843 844

845 846

847

848

849 850

851

852

853

854

855

856

857 858

859

860

861

862

863

864

865

867

868

869

870

871

872

873

- Zheng Xin Yong, Ruochen Zhang, Jessica Forde, Skyler Wang, Arjun Subramonian, Holy Lovenia, Samuel Cahyawijaya, Genta Winata, Lintang Sutawika, Jan Christian Blaise Cruz, Yin Lin Tan, Long Phan, Long Phan, Rowena Garcia, Thamar Solorio, and Alham Fikri Aji. 2023. Prompting multilingual large language models to generate code-mixed texts: The case of south East Asian languages. In *Proceedings* of the 6th Workshop on Computational Approaches to Linguistic Code-Switching, pages 43–63, Singapore. Association for Computational Linguistics.
- Ruochen Zhang, Samuel Cahyawijaya, Jan Christian Blaise Cruz, Genta Winata, and Alham Fikri Aji. 2023a. Multilingual large language models are not (yet) code-switchers. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12567–12582, Singapore. Association for Computational Linguistics.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. 2023b. Siren's song in the ai ocean: a survey on hallucination in large language models. *arXiv preprint arXiv:2309.01219*.
- Zhihan Zhang, Dong-Ho Lee, Yuwei Fang, Wenhao Yu, Mengzhao Jia, Meng Jiang, and Francesco Barbieri.
 2024. PLUG: Leveraging pivot language in crosslingual instruction tuning. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7025– 7046, Bangkok, Thailand. Association for Computational Linguistics.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Preprint*, arXiv:2306.05685.

A Dataset Details

A.1 Details of Source Data

CLIcK (Kim et al., 2024b)³ consists of 1,995 multiple-choice QA pairs, classified in two main categories (Culture, Language) and 11 sub-categories. CLIcK is sourced from various official Korean exams and textbooks, *e.g.*, College Scholastic Ability Test of Korea (CSAT). In this work, we only utilize data of eight sub-categories from Korean Culture category as our work aims to evaluate the effect of CS on activating Korean-specific knowledge.

³https://huggingface.co/datasets/EunsuKim/ CLIcK HAE-RAE (Son et al., 2024)⁴ is a Korean benchmark dataset originally crafted to capture cultural and contextual nuances inherent to the Korean language. We use 1,027 multiple-choice QA pairs regarding Korean culture. Both datasets are sourced from official Korean exams, textbooks, and text on the internet.

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

We combine two datasets and merge common categories (*i.e.*, Society, Geography, and Law), resulting in 2,372 QA pairs in nine categories: Popular, Economy, Politics, Tradition, General Knowledge, Society, Geography, History, and Law.

A.2 Dataset Statistics and License

We provide statistics of EnKoQA per domain in Table 2. We plan to release the dataset in public, under CC BY-NC license. We clarify that the source datasets are either open-source or used under authors' permission, ensuring that there are no issues regarding their use.

Domain	#
Economy	59
General	176
Geography	281
History	468
Law	435
Politics	84
Popular	41
Society	606
Tradition	222
Total	2,372

Table 2: Number of samples in EnKoQA.

A.3 Code-Mixing Index

We report CMI scores for our dataset in Table 3. In specific, we tokenized the sentence using bert-base-multilingual-cased, then removed all noisy tokens such as numbers or tags and counted the ratio of $\frac{num of Korean tokens}{num of all tokens}$. We report the distribution of QA accuracy on different CMI scores in Tradition and History, two domains where CS proved its effectiveness. If CMI is close to 0, sentence is mostly written in English, and close to 100 means vice versa. The number of samples at each end (0-10, 90-100) was very small, causing outliers. We can see that accuracy scores are

⁴https://huggingface.co/datasets/HAERAE-HUB/ HAE_RAE_BENCH_1.1

CMI		Tra	adition		History						
	Solar	Gemma 2 9B	Gemma 2 27B	GPT-40	Solar	Gemma2 9B	Gemma2 27B	GPT-40			
0–10	00.00	00.00	00.00	00.00	00.00	00.00	00.00	00.00			
10-20	100.0	100.0	100.0	100.0	50.00	50.00	100.0	50.00			
20-30	56.25	65.62	53.12	71.88	50.00	31.58	36.84	60.53			
30–40	72.34	59.57	55.32	76.6	66.67	48.72	45.3	77.78			
40–50	80.39	62.75	70.59	84.31	69.54	42.38	42.38	78.15			
50–60	83.72	74.42	74.42	93.02	62.35	40.00	35.29	74.12			
60–70	85.19	66.67	66.67	92.59	50.00	28.85	46.15	63.46			
70–80	66.67	58.33	66.67	91.67	57.89	15.79	21.05	57.89			
80–90	66.67	83.33	83.33	100.0	50.00	50.00	25.00	100.0			
90–100	00.00	00.00	00.00	00.00	00.00	00.00	00.00	00.00			

Table 3: Distribution of QA accuracy on different CMI scores in Tradition and History. If CMI is close to 0, sentence is mostly written in English, and close to 100 means vice versa. The number of samples at each end (0-10, 90-100) was very small, causing outliers.

quite evenly distributed across all ratios, suggesting that there is no distinct tendency between CMI and accuracy performance.

906 907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

It is important to note, however, that code switching metrics such as CMI, while offering a quantitative measure of token-level composition, are inherently limited in capturing the nuanced semantic and syntactic characteristics of code-switched texts. These metrics primarily rely on surface-level token ratios, which can inadvertently assign high scores to linguistically or contextually meaningless sequences. Consequently, they may over-represent the presence of meaningful code-switching patterns while failing to account for the deeper linguistic interplay that defines effective code-switching. For a more comprehensive discussion of these limitations, please refer to Srivastava and Singh, 2021.

A.4 Quality Control Guideline

We provide a guideline we used to filter the candidates and select the final candidate.

- Is the question written in English-Korean codeswitching, where matrix language is English and semantically important Korean words are embedded into English sentence?
- Do choices also follow the code-switched pattern of query?
- Does the syntactic structure of the sentence follow that of English?

• Are semantically important nouns and noun phrases from Korean sentence, and are they embedded into English sentence?

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

• Are functional words and grammatical morphemes kept in English?

A.5 Annotation Details

For dataset construction, two Korean native annotators with expert knowledge in Korean culture and equivalently fluent in English manually examine the candidates and select the most naturally codeswitched question, then cross-checked each other's assigned share of dataset. If a selected candidate appeared to be incorrect or suboptimal, the annotators engaged in thorough discussions until they reached an agreement on the most appropriate candidate.

Regarding inter-annotator agreement (IAA), although we did not compute a formal IAA score, significant effort was devoted to ensuring high annotation quality through extensive discussion and collaboration among annotators. In specific, the annotation process involved annotators who are fluent in both English and Korean are assigned each portion of the dataset to select a candidate for codeswitched question. Following this initial annotation, the annotators cross-checked each other's work to identify any discrepancies. If a selected candidate appeared to be incorrect or suboptimal, the annotators engaged in thorough discussions until they reached an agreement on the most appropriate candidate. This iterative and collaborative process was integral to constructing a high-quality dataset.

A.6 Dataset Size and Quality

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

985

986

987

988

989

990

991

992 993

994

995

996

997

998

999

Discussion on Dataset Size While we acknowledge the relatively limited size of EnKoQA dataset, we emphasize that quality often matters more than quantity as many studies (Pacchiardi et al., 2024; Maia Polo et al., 2024; Vivek et al., 2024) have demonstrated. Please note that we prioritized creating a high-quality dataset with rigorous manual validation and linguistic alignment, ensuring that the dataset serves as a reliable resource for codeswitching research. Additionally, while the size of Korean datasets is often limited given that Korean is a low-resource language, EnKoQA dataset is comparatively larger than the sizes of other Korean datasets. For instance, datasets in the Open Ko-LLM leaderboard (Park et al., 2024), such as Ko-ARC (1.1k), Ko-TruthfulQA (0.8k), and Ko-CommonGen (0.8k), are all smaller in scale than EnKoQA's 2,372 question-answer pairs. This highlights our effort to provide a relatively extensive resource within the constraints of dataset availability for minor languages.

Specifically, our quality control process includes human annotators thoroughly reviewing all LLMgenerated samples to assess the quality and naturalness. When any errors or unnatural code-switching patterns were identified, annotators corrected them to ensure that the final dataset adheres to high standards of our quality control. In that sense, GPT-3.5turbo served as an assistive tool for providing initial candidates, rather than generating final outputs. Therefore, we assert that any potential shortcomings of the translation tool were effectively mitigated through this meticulous human review and correction process.

Translating with GPT-3.5 We have conducted 1000 experiments on both GPT-3.5 and GPT-40 for trans-1001 lation and code-switching generation tasks. Inter-1002 estingly, we observed that after manual examina-1003 tion and correction process, the results from both 1004 models were comparable in terms of quality and 1005 naturalness. This is due to our rigorous human-in-1006 the-loop workflow that ensures any errors or un-1007 natural expressions are taken care of, regardless of 1008 the initial model used. Given this finding, we used 1009 GPT-3.5 for its cost efficiency while maintaining 1010 1011 high-quality standards through meticulous human examination and refinement. By prioritizing manual 1012 validation, we ensured that the final dataset reflects 1013 linguistic accuracy and naturalness, independent of 1014 the model used for preliminary generation. 1015

A.7 Data Sample

We also provide a sample of original Korean, trans-
lated English, and synthesized CS example question1017in Table 4. Note that unique terms or semantically
important words are properly embedded in Korean.1018

1016

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

B Experimental Details

B.1 Computational Resources and API Cost.

Llama3 and Gemma2 models. We used Huggingface model cards and run them on two NVIDIA A100 GPUs. Specifically, we used meta-llama/meta-llama-3-8b-instruct, meta-llama/meta-llama-3-70b-instruct, google/gemma-2-9b-it, google/gemma-2-27b-it.

GPT-3.5 and GPT-40. We used up-to-date versions of gpt-3-5-turbo and gpt-40 APIs. The cost for gpt-3-5-turbo was \$15 for EnKoQA generation and \$6 for experiment inference, while the cost for gpt-40 was \$23 for experiment inference.

Claude 3.5. We used claude-3-5-sonnet API from Anthropic AI^5 . The cost for claude-3-5-sonnet was \$21 for experiment inference.

Solar We used solar-mini API from Upstage⁶.

B.2 Prompts

We provide the following prompts used in our experiments. Table 12 contains the prompt used for generating code-switched text candidates across different levels of linguistic complexity. For QA inference tasks, we used the prompt presented in Table 13. The prompt for identifying relevant knowledge in a given context is provided in Table 14, while Table 15 shows the prompt used for leveraging this identified knowledge in downstream tasks.

B.3 Comparing CS and original Korean

We also present experimental results on original 1051 Korean data in Table 9. Generally, performances 1052 in original Korean are higher than in CS, while 1053 CS score approximates or equal to in many cases. 1054 Considering the English dominance of the LLMs, 1055 GPT-40 and Claude 3.5 Sonnet present advanced 1056 multilingual ability with over 80% accuracy. As 1057 our main research focus was towards on examining 1058 the effect of code-switching compared to dominant 1059

⁵https://www.anthropic.com/

⁶https://www.upstage.ai/

Lang	QUESTION	CHOICES
KO	다음 글의 (가)에 대한 (나)의 상대적 특성으로 옳은 것 은? (단, (가), (나)는 각각 겨울과 여름 중 하나임.) 우리나라는 더위와 추위에 대비하여 대청마루와 온돌 같은 전통 가옥 시설이 발달하였다. 대청마루는 바람 을 잘 통하게 하여 (가) 을 시원하게 지낼 수 있도록 설 치되었다. 온돌은 아궁이의 열을 방으로 전달하여 (나) 을 따뜻하게 지낼 수 있도록 설치되었다. 대청마루 는 중부와 남부 지역에 발달한 한편, 온돌은 대부분의 지 역에 발달하였다.	 (1) 평균 상대 습도가 높다. (2) 정오의 태양 고도가 높다. (3) 한파의 발생 일수가 많다. (4) 대류성 강수가 자주 발생한다. (5) 열대 저기압의 통과 횟수가 많다.
EN	What is the correct relative characteristic of (\downarrow^{+}) in re- lation to $(7\uparrow)$ in the following passage? (Note that $(7\uparrow)$) and (\downarrow^{+}) refer to either winter or summer.) In Korea, traditional house facilities such as daecheong- maru and ondol have developed to cope with heat and cold. Daecheongmaru is designed to allow good ventila- tion to keep $(7\uparrow)$ cool. Ondol transfers heat from the kitchen stove to the room to keep (\downarrow^{+}) warm. While daecheongmaru is developed in the central and southern regions, ondol is developed in most areas.	 The average relative humidity is high. The midday sun's altitude is high. There are many days of occurrence of cold waves. Heavy rainfall often occurs in Daeryuseong. There are many occurrences of passage of tropical cyclones.
CS	What is the correct relative characteristic of (나) in re- lation to (가) in the following passage? (Note that (가) and (나) refer to either winter or summer.) In 한국, 전통 가옥 시설 such as 대청마루 and 온돌 have developed to cope with heat and cold. 대청마루 is de- signed to allow good ventilation to keep (가) cool. 온돌 transfers heat from the kitchen stove to the room to keep (나) warm. While 대청마루 is developed in the 중부 and 남부 지역, 온돌 is developed in most areas.	 (1) The average 상대 습도 is high. (2) The 정오의 태양 고도 is high. (3) There are many days of occurrence of 한파. (4) 대류성 강수 often occurs. (5) There are many occurrences of passage of 열대 저기압.

Table 4: An example of Korean, English, and CS from dataset.

1060language, we exclude original Korean as our major1061concern. Also, there is a severe possibility that the1062performances may be influenced by already having1063seen the datasets, *i.e.*, data contamination, as the1064original datasets are sourced from official exams1065and texts from the Internet that are openly available.

1066 B.4 Open-ended QA.

1067 1068

1069

1070

1071

1072

1073

1074 1075

1076

1077

1078

1079

Out dataset, ENKOQA is multiple-choice QA dataset, following its original source datasets. We additionally explore the potential of code-switching on open-ended QA as well.

Results are shown in Table 5. Using same questions in our dataset, we instruct the model to respond in short answer and compute exact match score. It is noticeable that the performances are very low compared to multiple-choice QA results. We attribute this to the free-form response of open-ended tasks, causing more errors and hallucinations. It is observable that the models barely answer correctly in History and Popular.

C Evaluation Details	1080
C.1 Evaluation Criteria	1081
We provide evaluation guideline for human evalua- tion.	1082 1083
Faithfulness. Faithfulness evaluates the factual correctness of the knowledge.	1084 1085
• Knowledge list is very faithful. Every knowl- edge is factually correct.	1086 1087
• Knowledge list is somewhat faithful. Some, not every, knowledge is factually correct.	1088 1089
• Knowledge list is not faithful at all. Every knowledge is hallucinated.	1090 1091
Helpfulness. Helpfulness evaluates how useful the knowledge is for answering the question.	1092 1093
• Knowledge list is very helpful. Every knowl- edge is relevant to the question, and used for finding the answer.	1094 1095 1096
• Knowledge list is somewhat helpful. Some, not every, knowledge is useful for finding the answer.	1097 1098 1099

Model		Economy	Geography	History	Law	Politics	Popular	Society	Tradition
GPT-40 GPT-3.5	CS EN	85.00 80.00	20.00 00.00	40.00 05.00	30.00 05.00	30.00 10.00 75.00	05.00 00.00 45.00	50.00 05.00	35.00 00.00
	CS EN KOR	70.00 75.00 65.00	00.00 00.00 45.00	00.00 00.00 05.00	20.00 10.00 30.00	10.00 15.00 60.00	43.00 05.00 0.00 20.00	10.00 05.00 65.00	10.00 00.00 60.00
Llama3-70B	CS EN KOR	20.00 30.00 60.00	00.00 05.00 50.00	00.00 00.00 00.00	10.00 10.00 40.00	15.00 20.00 70.00	10.00 05.00 35.00	10.00 10.00 55.00	00.00 00.00 60.00
Llama3-8B	CS EN KOR	20.00 15.00 25.00	00.00 00.00 30.00	00.00 00.00 05.00	05.00 05.00 05.00	25.00 15.00 50.00	00.00 00.00 05.00	05.00 00.00 10.00	00.00 00.00 20.00

Table 5: QA performances on open-end QA.

• Knowledge list is not helpful at all. All knowledge are irrelevant with the question.

Pair-wise comparison. We comprehensively evaluate the quality of knowledge generated from CS and English questions in terms of both faithfulness and helpfulness. If both are identical, evaluators can choose Tie.

In case of LLM-as-a-judge evaluation, same criteria and instructions are given as prompts.

1109 C.2 Human Evaluator Qualifications

1100

1101

1102

1103

1104

1105

1106

1107

1108

For knowledge identification evaluation, collecting 1110 qualified bilingual evaluators was not easy due to 1111 the inherent challenge in code-switching research 1112 of necessitating fluent bilinguals as evaluators. Our 1113 dataset is composed of questions from Korean pro-1114 ficiency tests for foreigners and the Korean Col-1115 lege Scholastic Ability Test. Thus, it is designed 1116 at a level that would not be challenging for eval-1117 uators whom were born and raised in Korea, re-1118 ceived a Korean public education, and graduated 1119 1120 prestigious universities. We managed to collect four Korean graduate school students as our evaluators, 1121 all of whom are native Korean with sufficient un-1122 derstanding of Korean culture. Also, they possess 1123 qualified English exam scores, indicating that they 1124 have no problem in understanding Korean-English 1125 code-switched texts. To mitigate the shortage of 1126 labor force, we designed the evaluation criteria ob-1127 jectively, allowing for an assessment that is not 1128 subjective and has clear correct answers. Specif-1129 1130 ically, we evaluate knowledge identification based on two criteria: faithfulness and helpfulness. Faith-1131 fulness evaluates the factualness of the knowledge, 1132 so the evaluators are required to use their back-1133 ground knowledge as well as searching from faithful 1134

Model	History	Tradition	General	Law
GPT-40	0.41	0.64	0.62	0.62
Solar	0.26	-0.09	0.38	0.02
Gemma2 27B	0.25	0.52	0.17	0.34
Gemma2 9B	0.20	-0.07	0.05	0.24

Table 6: Cohen's kappa (κ) correlation scores between human and LLM-as-a-judge evaluation. Gray indicates poor agreement.

sources where gold knowledge exists. To evaluate 1135 helpfulness, evaluators are given a gold answer to 1136 the question and determine whether the knowledge 1137 is helpful for finding the answer, using their logical 1138 reasoning.

1139

1140

1143

1151

1152

1153

1154

1155

D Observations

In this section, we provide additional results and 1141 comprehensive observations throughout our work. 1142

D.1 Knowledge Identification Results

We observed that the majority of models benefitted 1144 from CS questions. Table 1 shows that scores in 1145 CS are higher on all models in Politics, and in case 1146 of Law, only three models (GPT-3.5, Llama3 70B, 1147 and Gemma2 9B) out of eight models performed 1148 worse. We can see in Average score, all models 1149 except Gemma2 27B performed better on CS. 1150

D.2 Knowledge Leveraging Results

We provide accuracy results of Knowledge Leveraging in Table 7. Figure 5 is a visualization of this table.

D.3 Error Analysis

We provide full results of error counts in Table 8. 1156 Note that as models get smaller and show poor per-1157

Model		Economy	General	Geography	History	Law	Politics	Popular	Society	Tradition	Average
GPT-40	CS EN	93.22 79.66	80.11 76.14	69.75 60.14	76.50 64.96	49.66 51.49	92.86 85.71	97.56 92.68	65.51 58.42	81.98 73.87	78.57 71.45
GPT-3.5	CS EN	74.58 69.49	37.50 49.43	39.15 43.06	30.13 34.62	32.41 34.02	82.14 73.81	75.61 65.85	50.50 47.03	63.06 55.41	53.90 52.52
Claude 3.5	CS EN	96.61 89.83	78.41 77.84	78.29 72.60	76.50 67.52	57.24 53.79	84.52 89.29	92.68 92.68	70.46 62.38	86.04 81.53	80.08 76.38
Solar	CS EN	83.05 88.14	53.98 53.98	52.31 47.69	62.61 37.61	40.46 39.08	85.71 76.19	78.05 70.73	55.78 51.65	72.97 65.32	64.99 58.93
Llama3 70B	CS EN	76.27 79.66	60.80 61.36	54.45 55.52	48.29 47.65	40.46 39.77	86.90 80.95	82.93 75.61	55.94 56.11	71.17 68.47	64.13 62.79
Llama3 8B	CS EN	76.27 72.88	39.20 37.50	40.57 40.57	38.46 32.91	33.33 31.72	72.62 72.62	73.17 75.61	47.03 50.33	56.31 59.01	53.00 52.57
Gemma2 27B	CS EN	77.97 79.66	51.70 55.68	48.04 50.18	41.24 36.11	36.78 40.46	78.57 69.05	78.05 75.61	53.63 52.15	65.32 61.26	59.03 57.80
Gemma2 9B	CS EN	76.27 67.80	50.57 44.32	44.84 45.20	40.60 35.68	35.63 39.08	77.38 70.24	73.17 65.85	53.14 49.50	61.71 61.71	57.03 53.26

Table 7: Knowledge leveraging performances of multilingual LLMs on CS and English settings. **Bold** indicates higher score between CS and English on each model. **Green** indicates the highest score from each domain.

Model		Economy	General	Geography	History	Law	Politics	Popular	Society	Tradition	Total
GPT-40	CS	0/0	0/0	0/0	0/0	0/2	0/0	0/0	0/5	0/1	0/8
	EN	0/8	0/1	1/11	0/15	2/15	1/3	0/0	1/40	0/6	5/99
GPT-3.5	CS	0/0	1/1	1/1	2/2	0/0	0/0	0/0	1/1	0/0	5/5
	EN	0/0	3/3	1/1	2/2	1/1	0/0	0/0	6/6	0/0	13/13
Claude 3.5 Sonnet	CS	0/0	0/0	0/0	0/0	0/0	0/0	0/0	0/0	0/0	0/0
	EN	0/0	0/0	0/0	0/0	0/1	0/0	0/0	0/4	0/0	0/5
Solar	CS	0/0	5/5	1/1	2/2	11/11	0/0	0/0	5/5	1/1	20/20
	EN	0/0	9/9	8/8	11/11	10/10	0/0	1/1	4/4	2/2	35/35
Llama3 70B	CS	2/2	6/6	6/6	11/11	3/3	0/0	1/1	3/3	2/2	34/34
	EN	2/2	7/7	10/10	24/24	12/12	2/2	0/0	4/4	2/2	63/63
Llama3 8B	CS	1/1	5/5	3/3	9/9	2/2	1/1	1/1	0/0	1/1	23/23
	EN	0/0	4/4	6/6	8/8	5/5	2/2	0/0	2/2	2/2	29/29
Gemma2 27B	CS	4/4	18/18	7/7	26/28	22/38	2/2	2/3	6/9	10/14	97/123
	EN	3/3	28/28	5/5	37/38	7/15	7/7	1/1	13/20	5/5	106/122
Gemma2 9B	CS	3/3	7/7	12/12	25/25	18/19	1/1	2/2	7/7	2/2	77/78
	EN	9/9	30/30	13/13	35/35	11/11	5/5	0/0	23/28	9/9	135/140

Table 8: Counts of None errors. Each cell indicates # of None errors / # of errors due to knowledge hallucination. **Bold** indicates that all errors are caused by hallucination.

formance in Korean, the number of errors increase. (See Gemma2 families.)

D.4 Case Study

1158

1159

1160

We examine a sample case to compare the capability 1161 of code-switching and English on knowledge acti-1162 vation. Table 10 shows the knowledge and answer 1163 generated by Solar in Tradition. The question asks 1164 about "정월대보름", a Korean traditional holiday 1165 that celebrates the first full moon of lunar new year. 1166 We observe that CS question preserves unique terms 1167 such as "정월대보름" and "귀밝이술" in Korean; 1168 1169 this helps the model to successfully activate faithful knowledge, consequently leading to the correct 1170 answer. However, in the case of English, not only 1171 are these cultural nuances lost in English question, 1172 but the model misunderstood the question to asking 1173

about "Dan-o", another Korean traditional holiday. Solar lacks in knowledge about "정월대보름" in English, or fails to activate encoded knowledge with its English translation.

1174

1175

1176

1177

We also provide a case of CS failing in knowl-1178 edge activation in Table 11. In the case of Gemma2 1179 9B on Law domain, hallucinations are observed 1180 in the knowledge generated from CS question. Ac-1181 cording to the Civil Act of the Republic of Korea, 1182 individuals under the age of 14 can only enter into 1183 binding contracts with the consent of their legal 1184 guardians. Additionally, individuals between the 1185 ages of 14 and 19 are not deprived of contractual 1186 effect; rather, they are granted the right to cancel 1187 such agreements at their discretion. Moreover, the 1188 knowledge generated in English incorrectly applies 1189

the U.S. standard, which defines minors as those
under 18 years of age, instead of the Korean stan-
dard, which applies to individuals under 19 years of
age. This finding suggests that English question is
not helpful for identifying necessary and language-
specific knowledge.



Figure 6: LLM-as-a-judge evaluation results on faithfulness between knowledge lists identified from CS and English questions.



Figure 7: LLM-as-a-judge evaluation results on helpfulness between knowledge lists identified from CS and English questions.



Figure 8: LLM-as-a-judge evaluation results on pairwise comparison between knowledge lists identified from CS and English questions.

Model		Economy	General	Geography	History	Law	Politics	Popular	Society	Tradition	Total
GPT-40	CS KO _{og}	91.53 94.92	78.41 76.70	69.04 75.09	74.79 76.50	55.86 58.62	90.48 89.29	95.12 97.56	63.70 67.00	85.14 85.59	78.23 80.14
GPT-3.5	CS KO _{og}	$\frac{71.19}{71.19}$	47.73 45.45	44.48 37.37	32.91 29.49	35.40 30.11	$\frac{70.24}{70.24}$	80.49 60.98	49.17 47.36	57.21 54.05	54.31 49.58
Claude 3.5 Sonnet	CS KO _{og}	93.22 89.83	72.16 71.59	72.95 80.78	73.08 76.28	62.53 66.67	86.90 89.29	<u>95.12</u> <u>95.12</u>	67.66 71.45	84.23 87.39	78.65 80.93
Solar	CS KO _{og}	83.05 84.75	55.11 51.70	54.09 55.87	63.46 64.74	42.76 43.22	80.95 82.14	85.37 82.93	54.29 55.28	75.23 76.58	66.03 66.36
Llama3 70B	CS KO _{og}	79.66 86.44	51.70 52.27	50.53 53.38	49.36 51.50	44.14 41.84	80.95 77.38	75.61 82.93	57.43 59.41	65.77 68.92	61.68 63.79
Llama3 8B	CS KO _{og}	69.49 72.88	40.34 38.07	36.30 37.37	35.68 36.75	35.63 35.40	75.00 72.62	73.17 70.73	45.05 51.32	54.05 55.86	51.63 52.33
Gemma2 27B	CS KO _{og}	79.66 83.05	46.02 45.45	$\frac{48.75}{48.75}$	41.03 45.94	45.29 41.38	$\frac{77.38}{77.38}$	78.05 80.49	54.79 56.60	65.32 67.12	59.59 60.68
Gemma2 9B	CS KO _{og}	79.66 77.97	42.05 41.48	44.13 45.91	40.17 38.89	41.15 40.92	73.81 80.95	$\frac{80.49}{80.49}$	53.30 53.80	65.77 66.22	57.84 58.51

Table 9: QA performances of LLMs on CS and original Korean dataset. CS scores are identical with Table 1. **Bold** indicates higher score, and <u>underline</u> indicates tie.

Example of Solar from the Tradition domain

[QUESTION]

The following is a question about 한국의 문화 and 전통.

Foreigners want to experience 한국의 전통 풍습 by celebrating the 정월 대보름. Which of the following is not part of the 체험내용?

(The following is a question about Korean culture and traditions.

Foreigners want to experience Korean traditional customs by celebrating the first full moon of the lunar year. Which of the following is not part of the experience?)

(1) Eat 오곡밥 made from five types of grains. (Eat five-grain rice made from five types of grains.)

(2) Drink 귀밝이술, wishing for better hearing. (Drink ear-brightening potion, wishing for better hearing.)

(3) Eat food made with 쑥 and 익모초 to overcome the heat. (Eat food made with mugwort and dried hydrangea to overcome the heat.)

(4) Munch on 견과류 like 땅콩, 호두, and 밤, hoping for no skin disease. (Munch on nuts like peanuts, walnuts, and chestnuts, hoping for no skin disease.)

[Knowledge from CS]

1. 정월 대보름 (Jeongwol Daeboreum): It is a traditional Korean holiday that takes place on the first full moon of the lunar new year. It is a time for families to gather, participate in various rituals, and enjoy traditional foods. 2. 오곡밥 (Ogokbap): A traditional rice cake made with five grains (rice, barley, soybeans, sesame seeds, and black beans) that is eaten during 정월 대보름 for good health and longevity.

3. 귀밝이술 (Gwi-bal-i-sul): A traditional Korean drink made from makgeolli (rice wine) that is consumed on 정월 대보름 to improve hearing and communication skills.

4. 쑥과 익모초 (Perilla and Mugwort): These herbs are used in traditional Korean medicine and are believed to have healing properties. While they may be used in some traditional Korean dishes, they are not specifically associated with 정월 대보름.

5. 견과류 (Nuts): Eating nuts like 땅콩 (peanuts), 호두 (walnuts), and 밤 (chestnuts) is a common practice during 정월 대보름 as they are believed to bring good luck and prevent skin diseases.

[ANSWER] (3)

[Knowledge from EN]

1. The first full moon of the lunar year is celebrated in Korea as a traditional holiday called "Dano."

2. One of the customs during Dano is eating "ongchimsik," a rice dish made with five grains.

3. Another custom is drinking "suriseong," a herbal drink believed to improve hearing.

4. Eating "yeot," a dish made with mugwort and dried hydrangea, is also a traditional practice during Dano to help overcome the heat.

5. Nuts like peanuts, walnuts, and chestnuts are often consumed during Dano as a symbol of good health and fortune.

6. Skin disease is not directly related to the traditional customs of Dano.

[ANSWER] (4) ...

Table 10: An example of knowledge activation by CS and English questions from the Tradition domain on Solar. The text in parentheses shows English translation.

Example of Gemma2 9B from the Law domain

[QUESTION]

What is the correct legal 판단 on the following case?

17세인 A, B, and C는 entered into a 계약 with 판매업자 D to purchase an expensive 노트북. A did not obtain 부모의 동의, while both B and C presented 부모의 동의서, but C forged the consent form. At the time of the contract, D did not know that A was a 미성년자. *Parents of A, B, and C are all 법정 대리인.

(What is the correct legal judgment on the following case?

17-year-olds A, B, and C entered into a contract with seller D to purchase an expensive laptop. A did not obtain parental consent, while both B and C presented parental consent forms, but C forged the consent form. At the time of the contract, D did not know that A was a minor. * Parents of A, B, and C are all legal representatives.)

(1) B's 부모 can cancel the 계약. (B's parents can cancel the contract.)

(2) The 부모 of A and C can both cancel the 계약. (The parents of A and C can both cancel the contract.)

(3) B cannot cancel the 계약 unlike A. (B cannot cancel the contract unlike A.)

(4) D cannot withdraw the 의사 표시 to conclude a 계약 with A. (D cannot withdraw the expression of intention to conclude a contract with A.)

(5) D can demand a prompt response from A regarding the 취소 of the 계약. (D can demand a prompt response from A regarding the cancellation of the contract.)

[Knowledge from CS]

1. **미성년자 계약의 효력**: 대한민국 민법상 14세 미만의 자는 계약력이 없으며, 14세 이상 19세 미만의 자는 부모의 동의가 없는 계약은 효력이 없다.

2. **부모의 동의**: 미성년자의 계약은 부모의 동의가 있어야 효력이 있다. 동의서가 위조된 경우, 계약은 무효이다.

3. **계약의 취소**: 계약의 당사자는 계약을 취소할 수 있다.

4. **의사표시의 취소**: 의사표시는 당사자의 의사에 따라 취소될 수 있다.

[ANSWER] (Correct answer is (3).) (2)

[Knowledge from EN]

1. In most jurisdictions, minors (typically under 18 years old) lack the legal capacity to enter into binding contracts.

2. Contracts entered into by minors are generally voidable at the option of the minor.

3. Parental consent can sometimes affect the validity of a contract entered into by a minor, depending on the specific laws of the jurisdiction.

4. Fraudulent misrepresentation, such as forging a parental consent form, can invalidate a contract.

[ANSWER] (2)

Table 11: An example of knowledge activation by CS and English questions from the Law domain on Gemma2 9B. The text in parentheses shows English translation.

Prompt for generating CS candidates

You are a bilingual who can speak both English and Korean fluently. I will give you a Korean and English (<Korean>, <English>) pair. They are semantically the same. Your job is to write Korean-English code-switching text with certain switching level by mixing Korean and English text. Example 1 is an example of code-switched text in different levels of code-switching. Use Example 1 as reference to understand the level of code-switching. Read the instructions carefully and solve the Task.

Instructions:

- Maintain English word order, that is, Subject-Verb-Object.

- Find semantically important given nouns and noun phrases from the text, and change {level} percent of them to Korean.

- Keep functional words in English.

- Keep the indicators such as (가), (나), ㄱ, ㄴ, 갑, 을 in Korean.

[Example 1]

<Korean>

제주도는 점성이 작고 유동성이 큰 마그마가 여러 차례 분출하여 형성된 방패 모양의 화산섬이다. 하지만 한라산의 정상부는 종 모양의 화산으로 이루어져 있으며, 산허리에는 오름으로 불리는 기생화산이 많이 형 성되어 있다.

<English>

Jeju Island is a shield-shaped volcanic island formed by multiple eruptions of small-sized and highly fluid magma. However, the top of Hallasan Mountain consists of a cone-shaped volcano, and many parasitic volcanoes called Oreum are formed on the hillsides.

<Code-switch with 30 percent of Korean>

Jeju Island is a shield-shaped 화산섬 formed by multiple eruptions of small-sized and highly fluid magma. However, the top of 한라산 consists of a cone-shaped volcano, and many 기생화산 called 오름 are formed on the hillsides.

<Code-switch with 50 percent of Korean>

Jeju Island is a 방패 모양의 화산섬 formed by multiple eruptions of small-sized and highly fluid 마그마. However, the top of 한라산 consists of a cone-shaped 화산, and many 기생화산 called 오름 are formed on the hillsides.

<Code-switch with 70 percent of Korean>

제주도 is a 방패 모양의 화산섬 formed by multiple eruptions of 크기가 작고 유동성이 큰 마그마. However, the top of 한라산 Mountain consists of a 종 모양의 화산, and many 기생화산 called 오름 are formed on the 산허리.

<Code-switch with 90 percent of Korean>

제주도 is a shield-shaped 화산섬 formed by multiple 분출 of small-sized and 유동성이 큰 마그마. However, the 정상부 of 한라산 consists of a cone-shaped 화산, and many 기생화산 called 오름 are formed on the 산허리.

[Task] <Korean> {question}

<English> {translation}

<Code-Switch>

Table 12: Prompt for generating code-switched text candidates in different levels.

Prompt for QA (CS)

You will be given a question and choices about Korea. The text are written in English-Korean code-switching, where matrix language is English and semantically important Korean words are embedded into English sentence. Your job is to answer the question. Read the [QUESTION] and choose the most appropriate answer from [CHOICES]. Only write your answer number in parentheses, like (1). Do not repeat the question or choice. Use Example 1 as a reference to answer Example 2.

<Example 1> [QUESTION] Which city is the 수도 of 한국?

[CHOICES] (1) 뉴욕 (New York) (2) 서울 (Seoul) (3) 파리 (Paris) (4) 도쿄 (Tokyo)

[ANSWER] (2)

<Example 2> [QUESTION] {question} [ANSWER]

Table 13: Prompt for QA inference.

Prompt for Knowledge Identification

You are a bilingual who is fluent in both Korean and English, and is knowledgeable about South Korea. You will be given a multiple choice question about South Korea. The text are written in English-Korean code-switching, where matrix language is English and semantically important Korean words are embedded into English sentence. Your job is to follow the instructions and write a list of knowledge that is necessary to know for solving the question correctly.

Instructions:

- Write a list of factual knowledge that are required for solving the question. Try to write each knowledge in one or two sentences. You can write in whichever language you can explain better, either Korean or English. Start this task with [KNOWLEDGE] tag.

- Only write knowledge that you definitely know. Do not write incorrect information.

- Do not repeat input text in your response. Do not generate new question. Stick to input text that is given to you.

I will give you an example for reference.

«Example 1»

[QUESTION]

Read the following question and choose the most appropriate answer. Who is the person who greatly defeated the soldiers of the 당나라 in the 안시성 싸움?

[CHOICES] (1) 양만춘 (2) 서희

(2) 서희 (3) 김유신 (4) 강감찬

(5) 윤관

[KNOWLEDGE]

1. 안시성 싸움 (Siege of Ansi): 안시성 싸움 (645 AD) was a famous military conflict between 고구려 and the 당 Dynasty. 고구려, under the leadership of 양만춘 (Yang Man-chun), successfully defended the 안시성 against the powerful 당 forces led by Emperor 태종.

2. 양만춘 (Yang Man-chun): He was the general who commanded the defense of 안시성, playing a key role in defeating the 당나라 army.

3. 서희 (Seo Hee): A 고려 diplomat famous for negotiating with the 거란 to avoid invasion, but not involved in the 안시성 싸움.

4. 김유신 (Kim Yu-shin): A general from the 신라 Kingdom, instrumental in the unification of the 한반도, but not involved in this specific battle.

5. 강감찬 (Gang Gam-chan): A 고려 military commander known for his victory over the 거란 in the 귀주대첩, unrelated to 안시성.

6. 윤관 (Yun Gwan): A 고려 general famous for his campaigns against the Jurchen, unrelated to the 한반도.

Now solve this. «Example 2» [QUESTION] {question}

[CHOICES] {choices}

Table 14: Prompt for Knowledge Identification task.

Prompt for Knowledge Leveraging

You are a bilingual who is fluent in both Korean and English, and is knowledgeable about South Korea. You will be given a multiple choice question and a list of knowledge that are relevant to the question. The text are written in English-Korean code-switching, where matrix language is English and semantically important Korean words are embedded into English sentence. Your job is to follow the instructions and select one choice from [CHOICES].

Instructions:

- Using given [KNOWLEDGE], explain concisely what and why you think is the answer. You can write in whichever language you can explain better, either Korean or English. Start this task with [EXPLANATION] tag. - Choose your final choice from [CHOICES]. The answer is one of the [CHOICES], so do not say 'none of the above'. You must write a index number in parentheses, like (1). Start this task with [ANSWER] tag.

- Do not repeat input text in your response. Do not generate new question. Stick to input text that is given to you.

I will give you an example for reference.

«Example 1»

[QUESTION]

Read the following question and choose the most appropriate answer. Who is the person who greatly defeated the soldiers of the 당나라 in the 안시성 싸움?

[CHOICES]

(1) 양만춘 (2) 서희 (3) 김유신 (4) 강감찬 (5) 윤관

[KNOWLEDGE]

1. 안시성 싸움 (Siege of Ansi): 안시성 싸움 (645 AD) was a famous military conflict between 고구려 and the 당 Dynasty. 고구려, under the leadership of 양만춘 (Yang Man-chun), successfully defended the 안시성 against the powerful 당 forces led by Emperor 태종.

2. 양만춘 (Yang Man-chun): He was the general who commanded the defense of 안시성, playing a key role in defeating the 당나라 army.

3. 서희 (Seo Hee): A 고려 diplomat famous for negotiating with the 거란 to avoid invasion, but not involved in the 안시성 싸움.

4. 김유신 (Kim Yu-shin): A general from the 신라 Kingdom, instrumental in the unification of the 한반도, but not involved in this specific battle.

5. 강감찬 (Gang Gam-chan): A 고려 military commander known for his victory over the 거란 in the 귀주대첩, unrelated to 안시성.

6. 윤관 (Yun Gwan): A 고려 general famous for his campaigns against the Jurchen, unrelated to the 한반도.

[EXPLANATION]

The question specifically asks about the 안시성 싸움 (Siege of Ansi) and who defeated the 당나라 soldiers in that battle. Based on historical facts, the leader who played a key role in defending 안시성 and defeating the 당나라 army was 양만춘 (Yang Man-chun).

[ANSWER] (1)

Now solve this. «Example 2» [QUESTION] {question} [CHOICES] {choices}

[KNOWLEDGE] {knowledge}

Table 15: Prompt for Knowledge Leveraging task.