## **Controlling Language Confusion in Multilingual LLMs**

#### **Anonymous ACL submission**

#### Abstract

Large language models often suffer from language confusion, a phenomenon where responses are partially or entirely generated in unintended languages. This can critically impact user experience in low-resource settings. We hypothesize that conventional supervised finetuning exacerbates this issue because the softmax objective focuses probability mass only 010 on the single correct token but does not explicitly penalize cross-lingual mixing. Interestingly, by observing loss trajectories during the pretraining phase, we observe that models fail to learn to distinguish between monolingual and language-confused text. Additionally, we find 016 that ORPO, which adds penalties for unwanted output styles to standard SFT, effectively sup-017 presses language-confused generations even at 018 high decoding temperatures without degrading overall model performance. Our findings suggest that incorporating appropriate penalty terms can mitigate language confusion in low-022 resource settings with limited data.

#### 1 Introduction

024

037

041

Scaling large language models has empirically delivered substantial gains in multilingual capabilities (Hurst et al., 2024; Cohere et al., 2025; Yang et al., 2025), among diverse tasks such as machine translation (Alves et al., 2024), summarization (Forde et al., 2024), and reasoning (Son et al., 2025). However, despite their growing capabilities, LLMs often suffer from language confusion (Marchisio et al., 2024), a failure mode in which outputs inadvertently blend multiple languages. This hampers real-world deployment of LLM systems as even the most minor language confusion may be critical to user experience (Son et al., 2024a). This issue is particularly pronounced in low-resource settings, where limited supervision exacerbates cross-lingual interference (Arivazhagan et al., 2019; Wang et al., 2023).

However, little research has been conducted on *why* such behavior may happen. In this work, we draw inspiration from the training methodology proposed by Hong et al. (2024), which applies supervised fine-tuning to preferred generation styles while imposing penalties on disfavored ones.

042

043

044

047

049

051

054

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

079

In this work, we conduct two experiments to investigate whether language confusion arises from the absence of an explicit penalty against undesired languages. First, we track the training loss of two model families (SmolLM2 (Allal et al., 2025) and OLMo2 (OLMo et al., 2024)) throughout their pretraining process. In both cases, the loss of language-confused outputs steadily decreases over time, indicating that the models do not learn to disfavour confused generations. Additionally, by using ORPO (Hong et al., 2024) for an additional three epochs of fine-tuning, we show that introducing an explicit penalty against unwanted languages effectively restricts cross-lingual mixing.

#### 2 Preliminaries

#### 2.1 Related Works

What is language confusion? Language confusion, also known as language mixing or codemixing, occurs when two or more languages are mixed within a single utterance (Chen et al., 2024; Yoo et al., 2024). This phenomenon is particularly prevalent in low-resource languages (Arivazhagan et al., 2019) and even appears in state-of-the-art models (u/VictorRM, 2025). Diverse discussions have emerged regarding language confusion. Although it can sometimes support multilingual transfer (Wang et al., 2025), mixed-language responses may undermine user experience, as they can be perceived as signs of incompetence (Son et al., 2024a).

#### 2.2 Quantifying Language Confusion

Measuring language confusion may be challenging as LLM-Judges (Zheng et al., 2023) remain unreliable (Son et al., 2024b), and rule-based methods
cannot distinguish genuine confusion from legitimate uses of foreign-language (e.g., abbreviations).
In this work, we leverage two metrics Word Precision Rate (WPR) and Language Precision Rate (LPR) proposed by Marchisio et al. (2024).

**WPR** computes the overall fraction of tokens produced in the target language, offering a granular view of how consistently a model sticks to one language. Where  $\mathcal{T} = \bigcup_{i=1}^{N} T_i$  is the set of all valid tokens across N outputs WPR is defined as:

091

100

101

102

103

104

106

107

108

110

111

112

113

114

115

116

117

118 119

120

121

122

123

$$\frac{\left|\left\{t \in \mathcal{T} : \text{is}_{\text{Korean}}(t)\right\}\right|}{|\mathcal{T}|} \tag{1}$$

**LPR** counts the proportion of sentences in which at least 90 % of tokens belong to the target language, thereby penalizing any cross-lingual intrusions. Given  $I(\cdot)$  is the indicator function and  $s_i$ the *i*-th sentence LPR is defined as:

$$\frac{1}{N} \sum_{i=1}^{N} I\left(\frac{\left|\{t \in s_i : \text{is}\_\text{Korean}(t)\}\right|}{\left|\{t \in s_i : \text{is}\_\text{valid}(t)\}\right|} \ge 0.9\right)$$
(2)

Additionally, as noted above, rule-based metrics alone cannot distinguish true language confusion from minor lexical variations, such as numerals, named entities, or common loanwords. Therefore, alongside WPR and LPR, we also report the proportion of responses with WPR and LPR exceeding 0.9. Empirically, we observe that many such responses remain perfectly acceptable sentences containing a few legitimate English terms. For examples with varying WPR and LPR levels of sentences, see Appendix D.

#### **3** Experimental Setup

#### **3.1 Dataset Preperation**

To facilitate pairwise preference learning, we constructed instruction-centered triplet datasets. Each triplet comprises a Korean prompt *(input)*, a fully Korean response *(chosen)*, and an alternative response exhibiting code-mixing or full unexpected language *(rejected)*.

We construct three multilingual datasets based on existing Korean corpora, each designed to represent a different form of language confusion. The OIG dataset (LAION, 2022; Heegyu, 2023) and HC3 dataset (Guo et al., 2023; Na, 2023) pair Korean prompts with rejected responses written entirely in English. In contrast, the KoAlpaca ### Input: 여행 준비를 위한 최고의 팁은 무엇입니까? ### Chosen: 1.간식과 물과 같은 물품이 충분한지 확인하십시오. 2.경로를 미리 계획하여 목적지와 도착하는 데 걸리는 시간을 알 수 있습니다. 3.짐은 가볍게 하되 재킷, 모자, 장갑 등을 준비하십시오. ### Rejected: 1.Make sure you have enough supplies, such as snacks and water. 2.Plan your route in advance so that you know where you're going and how long it will take to get there. 3.Pack light but still be prepared with items like jackets, hats, gloves, etc.

Figure 1: Dataset structure (OIG, Chosen-Rejected pair)

dataset (Beomi, 2023) introduces more nuanced confusion by synthetically injecting translated English or Chinese tokens into Korean outputs, resulting in code-mixed responses. Additional preprocessing and filtering steps are described in Appendix A. 124

125

126

127

128

129

130

131

132

133

134

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

160

#### 3.2 Experiment Setup

We fine-tuned two publicly available instructiontuned language models—SmolLM2-1.7B (Allal et al., 2025) and OLMo2-7B (OLMo et al., 2024), selected for their ability to generate Korean text among lightweight open source models. Detailed training configurations are provided in Appendix B.

#### **3.3 Evaluation Protocol**

We evaluate three model variants: **Base**, the original instruction-tuned model; **SFT**, supervised finetuned on Korean prompt–response pairs from the OIG dataset; and **ORPO**, fine-tuned using Odds Ratio Preference Optimization, on the same dataset.

#### 4 Main Results

Prior work shows LLMs default to high-frequency, dominant-language tokens when uncertain, causing language confusion (Marchisio et al., 2024). We hypothesize that the standard next-token prediction objective exacerbates this bias: softmax focuses probability mass on the correct token but does not explicitly penalize cross-lingual mixing.

# 4.1 Loss-Based Diagnostic: Do LLMs Penalize Language Mixing?

We begin with the observation that, during pretraining, neither SmolLM2 (Allal et al., 2025) model learns to penalize language confusion, as shown by their loss trajectories in Figure 2.

In principle, a model that internalizes a robust linguistic preference should learn to assign lower loss to coherent Korean-only generations while preserving relatively higher loss for language-confused

	Model			SmolL	M2-1.7B					OLMo2-7B			
	Temperature	0.7		1.0		1.2		0.7		1.0		1.2	
		Base	ORPO	Base	ORPO	Base	ORPO	Base	ORPO	Base	ORPO	Base	ORPO
Metric	WPR > 0.9 ratio LPR > 0.9 ratio Average WPR Average LPR	96.1% 92.6% 0.9821 0.9681	100.0% 99.9% 0.9999 0.9996	94.3% 88.5% 0.9696 0.9496	100.0% 100.0% 1.0 1.0	81.4% 71.2% 0.8953 0.8434	100.0% 99.9% 0.9999 0.9999	96.3% 71.2% 0.9818 0.9379	99.8% 99.7% 0.9998 0.9992	91.8% 46.0% 0.9576 0.8684	99.9% 99.8% 0.9998 0.9995	7.5% 0.5% 0.6799 0.3044	99.0% 96.8% 0.9962 0.9881

Table 1: Comparison of SmolLM2 and OLMo2 models across temperatures (Base vs. ORPO). All metrics are higher-is-better: higher values indicate stronger language consistency.



Figure 2: Average loss for monolingual and code-mixed responses across training tokens (SmolLM2)

outputs. Contrary to expectations, we observe a monotonic decrease in loss for both chosen and rejected responses. This trend may suggest that, in the absence of explicit preference signals, models eventually learn to prefer *any* sequence of tokens they have seen during training, without distinguishing linguistically coherent and code-mixed outputs. Such behavior persists up to 7B scale, suggesting that model size alone cannot resolve the issue. See Appendix C for results of OLMo2 models (OLMo et al., 2024).

161

162

163

164

165

166

167

168

170

171

172

173

174

175

176

177

178

179

181

183

184

187

#### 4.2 Generation-level evaluation: WPR and LPR Comparison

To evaluate the effectiveness of preference-based tuning method, we compare the generation performance of the Base and ORPO-tuned models using WPR and LPR under varying decoding temperatures. Each model generated responses for the same set of 1,000 prompts, repeated three times per prompt, and all reported scores are averaged across the three generations.

As summarized in Table 1, we observe the following trends:

• ORPO-tuned models consistently outperform the base-instruct models, achieving near-perfect WPR and LPR even at high temperature settings (up to 1.2). • Temperature significantly impacts the Base models, with LPR dropping below at temperature 1.2, indicating severe language confusion in the absence of preference-based finetuning.

### 5 Additional Results

## 5.1 Comparison with other fine-tuning methods

To assess how ORPO compares to other common fine-tuning paradigms, we conducted additional experiments using SFT under identical conditions. The results are summarized in Appendix E for SmolLM2 and OLMo2, respectively.

We observe that across both model families, ORPO consistently achieves strong WPR and LPR scores, matching or slightly outperforming SFT.

To further investigate whether preference-based learning offers additional internal modeling advantages, we conduct a loss-based diagnostic analysis on the evaluation subset HC3 and compare the loss between original (*chosen*) and code-mixed (*rejected*) responses.



Figure 3: Loss of SmolLM2 models across tuning methods for both original and code-mixed responses

We found that ORPO assigns significantly higher loss to code-mixed responses compared to other models, indicating stronger penalization of

208

209

210

211

212



Figure 4: Loss of OLMo2 models across tuning methods for both original and code-mixed responses

language-confused outputs. On the HC3 evaluation set, ORPO yields an average delta loss of 0.8379 for SmolLM2 and 4.6778 for OLMo2-both the highest among all fine-tuning methods. This increased separation suggests that ORPO fine-tuning more effectively reinforces internal preferences for linguistically consistent outputs, enabling more reliable discrimination between coherent and code-mixed generations (Figure 3, Figure 4).

213

214

215

216

217

218

219

221

222

224

226

231

232

237

240

241

243

244

247

This increased separation indicates that ORPO fine-tuning more effectively reinforces internal preferences for linguistically consistent outputs, allowing the model to more reliably discriminate between coherent and code-mixed generations.

# 5.2 Does ORPO Fine-Tuning Lead to a Trade-off in General QA Capabilities?

We assess whether ORPO fine-tuning, which mitigates language confusion, adversely affects general performance by evaluating our models on the HAE-RAE benchmark—a Korean multiple-choice QA suite covering general knowledge, history, loanwords, and rare vocabulary (Son et al., 2023). We omit more challenging reasoning benchmarks due to the modest size of our models and limited training data. We compared three models: Base, SFT, and ORPO fine-tuned model.

Figure 5 reports the average accuracies in all subcategories for the SmolLM2 and OLMo2 models. The results show no statistically significant performance degradation in the three tuning methods.

These findings suggest that neither SFT nor ORPO introduces measurable harm to general QA capabilities. In particular, ORPO maintains general QA performance while reducing language confusion.



Figure 5: Average accuracy across training methods for SmolLM2 and OLMo2.

248

249

250

251

252

254

255

256

257

258

259

260

261

262

263

264

265

#### 6 Conclusion

This work investigates why language confusion may happen in multilingual large language models and empirically show that penalizing unwanted languages via preference optimization is effective in restricting such behavior. By analyzing loss trajectories and generation-level evaluations, we demonstrate that preference-based tuning could enhance the model's ability to distinguish between original and language-confused outputs while preserving general question answering capabilities. These results suggest that incorporating explicit preference signals during fine-tuning provides a promising approach for reinforcing linguistic fidelity in multilingual settings. Moreover, we suggest that future research may explore the use of penalty terms even in the pretraining phase to penalize language confusion earlier in the training effectively.

#### References

- Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, and 1 others. 2025. Smollm2: When smol goes big–datacentric training of a small language model. arXiv preprint arXiv:2502.02737.
  - Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, and 1 others. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.
- Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Roee Aharoni, Melvin Johnson, and Wolfgang Macherey.
   2019. Massively multilingual neural machine translation in the wild: Findings and challenges. arXiv preprint arXiv:1907.05019.
- Beomi. 2023. Koalpaca: Korean instruction-tuning dataset.
- Yiyi Chen, Qiongxiu Li, Russa Biswas, and Johannes Bjerva. 2024. Large language models are easily confused: A quantitative metric, security implications and typological analysis. *arXiv preprint arXiv:2410.13237*.
- Team Cohere, Arash Ahmadian, Marwan Ahmed, Jay Alammar, Yazeed Alnumay, Sophia Althammer, Arkady Arkhangorodsky, Viraat Aryabumi, Dennis Aumiller, Raphaël Avalos, and 1 others. 2025. Command a: An enterprise-ready large language model. *arXiv preprint arXiv:2504.00698*.
- Jessica Zosa Forde, Ruochen Zhang, Lintang Sutawika, Alham Fikri Aji, Samuel Cahyawijaya, Genta Indra Winata, Minghao Wu, Carsten Eickhoff, Stella Biderman, and Ellie Pavlick. 2024. Re-evaluating evaluation for multilingual summarization. In *Proceedings* of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 19476–19493, Miami, Florida, USA. Association for Computational Linguistics.
- Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection. *arXiv preprint arXiv*:2301.07597.
- Heegyu. 2023. Oig-small-chip2-ko. https: //huggingface.co/datasets/heegyu/ OIG-small-chip2-ko.
  - Jiwoo Hong, Noah Lee, and James Thorne. 2024. Orpo: Monolithic preference optimization without reference model. *arXiv preprint arXiv:2403.07691*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1

others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.

- LAION. 2022. Open instruction generalist (oig) dataset. https://laion.ai/blog/oig-dataset/.
- Kelly Marchisio, Wei-Yin Ko, Alexandre Bérard, Théo Dehaze, and Sebastian Ruder. 2024. Understanding and mitigating language confusion in llms. *arXiv preprint arXiv:2406.20052*.
- Yohan Na. 2023. Hc3-ko: Korean human chatgpt comparison corpus. https://huggingface.co/ datasets/nayohan/HC3-ko. Accessed: 2025-05-17.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, and 1 others. 2024.
  2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728– 53741.
- Guijin Son, Jiwoo Hong, Hyunwoo Ko, and James Thorne. 2025. Linguistic generalizability of test-time scaling in mathematical reasoning. *arXiv preprint arXiv*:2502.17407.
- Guijin Son, Hyunwoo Ko, Hoyoung Lee, Yewon Kim, and Seunghyeok Hong. 2024a. Llm-as-a-judge & reward model: What they can and cannot do. *arXiv preprint arXiv:2409.11239*.
- Guijin Son, Hanwool Lee, Suwan Kim, Huiseo Kim, Jaecheol Lee, Je Won Yeom, Jihyu Jung, Jung Woo Kim, and Songseong Kim. 2023. Hae-rae bench: Evaluation of korean knowledge in language models. *arXiv preprint arXiv:2309.02706*.
- Guijin Son, Dongkeun Yoon, Juyoung Suk, Javier Aula-Blasco, Mano Aslan, Vu Trong Kim, Shayekh Bin Islam, Jaume Prats-Cristià, Lucía Tormo-Bañuelos, and Seungone Kim. 2024b. Mm-eval: A multilingual meta-evaluation benchmark for llm-as-a-judge and reward models. *arXiv preprint arXiv:2410.17578*.
- u/VictorRM. 2025. O3 thinks in chinese for no reason randomly. Reddit, r/OpenAI. Accessed 2025-05-19.
- Bin Wang, Zhengyuan Liu, Xin Huang, Fangkai Jiao, Yang Ding, AiTi Aw, and Nancy F Chen. 2023. Seaeval for multilingual foundation models: From cross-lingual alignment to cultural reasoning. *arXiv preprint arXiv:2309.04766*.
- Zhijun Wang, Jiahuan Li, Hao Zhou, Rongxiang Weng, Jingang Wang, Xin Huang, Xue Han, Junlan Feng, Chao Deng, and Shujian Huang. 2025. Investigating and scaling up code-switching for multilingual language model pre-training. arXiv preprint arXiv:2504.01801.

363

364

366

367

368

369

370

371

372

373

374

266

269

271

273

274

275

276

277

278

279

280

281

282

290

296

297

299

301

303

307

310

311

312

314

315

316

317

319

320

451

452

453

454

455

456

457

458

423

424

- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. arXiv preprint arXiv:2505.09388.
- Haneul Yoo, Cheonbok Park, Sangdoo Yun, Alice Oh, and Hwaran Lee. 2024. Code-switching curriculum learning for multilingual transfer in llms. *arXiv preprint arXiv:2411.02460*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.

#### A Dataset preprocessing

375

376

379

394

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

**KoAlpaca (Code-Mixed Rejection):** We constructed this dataset using the KoAlpaca<sup>1</sup> corpus, a Korean instruction-tuning dataset modeled after Stanford Alpaca (Beomi, 2023). Each triplet contains a Korean instruction, a fully Korean chosen response, and a synthetically generated code-mixed rejected response, created by injecting randomly selected English or Chinese tokens—translated via the Google Translate API—at random word-level positions.

To ensure high linguistic purity, we applied the following preprocessing steps: (1) filtered for chosen responses written entirely in Korean, guaranteeing a Word-level Pass Rate (WPR) and Line-level Precision Rate (LPR) of 1.0; (2) applied string normalization (e.g., whitespace trimming) to instruction, chosen, and rejected fields.

**OIG (Fully English Rejection):** We constructed a triplet dataset using the OIG-small-chip2-ko<sup>2</sup> corpus, which contains over 210K instructionresponse pairs translated into Korean from the original English OIG dataset (LAION, 2022). Each triplet comprises a Korean instruction, a fully Korean chosen response, and a fully English rejected response. This dataset is designed to evaluate the model's ability to distinguish between clearly separated linguistic domains.

We applied several preprocessing steps to improve data quality: (1) applied string normalization; (2) filtered for chosen responses containing only Korean text; (3) discarded samples where the length ratio between chosen and rejected responses fell outside the range of 0.4 to 2.0; (4) removed duplicate instructions. Each dataset contains approximately 10,000 instruction-response triplets, selected for linguistic consistency and diversity.

**HC3 (Fully English Rejection):** We also constructed quadruplet dataset using the HC3-ko<sup>3</sup>, which contains 24.3k insruction-response by human-response by GPT triplet translated into Korean from the original English HC3 dataset(Guo et al., 2023)(Na, 2023).

Each quadruplet comprises a Korean instruction, fully Korean chosen response, fully English rejected response, and a confusion response where eight Korean words are randomly selected and replaced with their translations into either English or Chinese. This dataset is designed to evaluate the model's generalizing ability to use the unseen data during training.

We applied several preprocessing steps to improve data quality: (1) applied string normalization; (2) filtered for chosen responses containing only Korean text; (3) discarded samples where the length ratio between chosen and rejected responses fell outside the range of 0.4 to 2.0; (4) removed duplicate instructions. (5) removed responses exhibiting generation failures caused by the language model, such as repeated phrases or malformed outputs due to server errors.

#### **B ORPO** Training Configuration

Table 2 outlines the training configuration used for ORPO fine-tuning. Both SmolLM2-1.7B and OLMo-2-1124-7B were trained for 3 epochs with a global batch size of 128. ORPO's weighting coefficient  $\beta$  was set to 0.1 across experiments, and training was performed using the DeepSpeed ZeRO-2 framework.

Parameter	SmolLM2-1.7B	OLMo-2-1124-7B
GPUs	A6000 × 1	H100 × 2
Max sequence length	8192	4096
Micro batch size	8	8
Gradient accumulation	16	8
Global batch size	128	128
Training steps	223	223
Epochs	3	3
ORPO $\beta$ value	0.1	0.1
Optimizer	AdamW	AdamW
Framework	DeepSpeed ZeRO-2	DeepSpeed ZeRO-2

Table 2: Training configuration for ORPO fine-tuning on SmolLM2 and OLMo2 models.

<sup>&</sup>lt;sup>1</sup>https://huggingface.co/datasets/beomi/ KoAlpaca-v1.1a

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/datasets/heegyu/ OIG-small-chip2-ko

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/datasets/nayohan/ HC3-ko



Figure 6: The average loss of original (monolingual) and code-mixed responses across training checkpoints for OLMo2 models.

460

471

472

### To assess whether the failure to penalize language confusion generalizes across architectures, we also tracked the loss trajectories of OLMo2 models (1B and 7B) throughout pretraining. As shown in Figure 6, both original and code-mixed responses exhibit a steady decrease in loss, mirroring the trend observed in SmolLM2 (Figure 2). Despite the increase in model capacity, the gap between two responses does not widen. This suggests that pretraining objectives alone may not induce meaningful linguistic preferences.

# D Samples of different levels of WPR and LPR

To enable interpretable comparisons across models, 473 we report the proportion of generations that exceed 474 a threshold of 0.9 for both WPR and LPR. This 475 threshold was chosen based on manual inspection 476 by a native Korean speaker (the author), who re-477 viewed a large number of generated samples and 478 heuristically identified 0.9 as a practical cutoff that 479 separates mostly monolingual responses from visi-480 bly code-mixed ones. This level of tolerance allows 481 minor lexical variation (e.g., loanwords, numerals) 482 while still maintaining strong target-language align-483 ment. It also aligns with real world expectations 484 485 for language consistency, particularly in Korean, where partial foreign-language inclusions are not 486 uncommon but still undesirable in many contexts. 487 Representative examples illustrating this threshold-488 ing effect are shown in Figure 7. 489

## E Generation-level evaluation: other models

In addition to ORPO, we evaluate two other finetuning methods: Supervised Fine-Tuning (SFT) and Direct Preference Optimization (DPO). 490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

508

509

510

511

512

513

514

515

516

517

518

519

520

Direct Preference Optimization (DPO) is a preference-based tuning method that trains models to maximize the log-probability margin between preferred and rejected responses (Rafailov et al., 2023).

Table 3 describes the detailed training configurations used for DPO fine-tuning. All settings were selected to closely match the original DPO implementation where possible.

Table 4 and Table 5 summarize the generation performance of SmolLM2 and OLMo2 models across three decoding temperatures (0.7, 1.0, 1.2) and three fine-tuning methods (SFT, DPO, ORPO). We report four key metrics: the ratio of outputs with WPR > 0.9, LPR > 0.9, average WPR, and average LPR.

Across both model families, ORPO consistently outperforms DPO and performs on par with or slightly better than SFT in terms of language fidelity. In particular, ORPO maintains near-perfect WPR and LPR values across all temperature settings, while DPO exhibits significant degradation at higher temperatures, most notably on the OLMo2 model at temperature 1.2 (LPR > 0.9 ratio drops to 52.1%. SFT remains relatively stable across temperatures.

Parameter	SmolLM2-1.7B (DPO)	OLMo-2-7B (DPO)
GPUs	A6000 × 1	A6000 × 4
Dataset size	10,000	10,000
Max sequence length	8192	4096
Micro batch size	8	4
Gradient accumulation	8	4
Global batch size	64	64
Training steps	467	467
DPO $\beta$ value	0.1	0.1
Optimizer	RMSprop	RMSprop
Framework	DeepSpeed ZeRO-2	DeepSpeed ZeRO-2

Table 3: DPO training configuration for SmolLM2 and OLMo2 models.

Metric	temperature = 0.7			tem	perature =	= 1.0	temperature = 1.2		
	SFT	DPO	ORPO	SFT	DPO	ORPO	SFT DPO		ORPO
WPR > 0.9 ratio LPR > 0.9 ratio Average WPR Average LPR	99.9% 99.8% 0.9998 0.9994	94.2% 92.3% 0.9760 0.9705	100.0% 99.9% 0.9999 0.9996	100.0% 100.0% 1.0000 1.0000	96.9% 94.4% 0.9857 0.9780	100.0% 100.0% 1.0000 1.0000	100.0% 99.7% 0.9998 0.9993	95.0% 90.5% 0.9823 0.9629	100.0% 99.9% 0.9999 0.9999

Table 4: Performance of SmolLM2 across temperature and tuning methods (SFT, DPO, ORPO)

Table 5: Performance of OLMo2 across temperature and tuning methods (SFT, DPO, ORPO)

Metric	temperature = 0.7			tem	perature :	= 1.0	temperature = 1.2		
	SFT	DPO	ORPO	SFT	DPO	ORPO	SFT	DPO	ORPO
WPR > 0.9 ratio	99.8%	99.5%	99.8%	99.9%	99.4%	99.9%	99.1%	94.4%	99.0%
LPR > 0.9 ratio	99.7%	92.7%	99.7%	99.8%	89.4%	99.8%	96.8%	52.1%	96.8%
Average WPR	0.9996	0.9959	0.9998	0.9998	0.9938	0.9998	0.9970	0.9649	0.9962
Average LPR	0.9988	0.9847	0.9992	0.9997	0.9791	0.9995	0.9915	0.8897	0.9881

Input	Generated	WPR	LPR
임신 중 언제부터 수영 을 할 수 있나요? 안녕 하세요. 제 가족 중에 임 신하신 분이 계시는데 수영을 하고자 하셔서 요. 몇 개월부터 수영이 가능한지 궁금합니다.	임신 1달부터 수영이 가능합니다. 하지만, 수영 중 피부 와 기관계의 건강을 위해, 수영 시간이 짧게 시작하고, 수 영 및 노출 시간을 조절하여 피부에 부담을 최소화하는 것이 좋습니다. 또한, 임신 중에 수영을 하는 것은 개인의 건강 상태, 임신의 진화 상태, 및 관련 건강 문제 등을 고 려해야 합니다.	1.0	1.0
개량 발명과 회피 발명 은 무엇인가요? 또한 이 러한 발명은 누구에 의 해 할 수 있을까요? 마 지막으로 원천 기술은 무엇인지 설명해주세요.	개량 발명(改良發明)은 기존 발명의 기능이나 성능을 향 상시키기 위해 수정·보완한 기술을 말합니다. 회피 발명 (回避發明)은 기존 특허를 침해하지 않도록 우회적인 방 법으로 개발한 기술이며, 연구개발자(R&D)나 발명가에 의해 이루어질 수 있습니다. 반면, 원천 기술(源泉技術)은 특정 기술 분야에서 새로운 paradigm을 제시하는 근본 적인 기반 기술로, 주로 연구소, 대학, 또는 대기업의 연 구진에 의해 개발됩니다.	0.9	0.67
건전지의 플러스 극과 마이너스 극의 역할은 무엇인가요?	건전지의 플러스 극(positive terminal)은 전자를 받는 역할을 하며, 이는 전기 회로에서 전류가 흐르는 방향의 starting point입니다. 반대로, 마이너스 극(negative terminal)은 전자를 내보내는 역할을 하며, 에너지가 회 로를 통해 흐르는 것을 가능하게 합니다. 이러한 양극 (anode)과 음극(cathode)의 역할은 전기적 potential 차이를 기반으로 작동합니다. 건전지를 사용할 때는 polarity를 정확히 확인하는 것이 중요합니다.	0.8	0.0
커피와 설탕과 프림을 섞었을 때 흔들면 커피 가 위로 올라오는 이유 는 무엇인가요?	가장 큰 이유는 설탕, 프림, 커피의 밀도 차이 때문입니다. 설탕과 프림은 커피보다 밀도가 높아 아래로 가라앉습니 다. 그 과정에서 커피는 상대적으로 위로 밀려 올라가게 됩니다. This phenomenon is caused by the difference in density among the components. As heavier particles sink, lighter coffee is displaced upward through convection-like motion.	0.5	0.6

Figure 7: Samples of generated responses at varying WPR and LPR levels