

# Towards Confident Multilingual Generation from English-Centric LLMs: A Tuning-Free Approach

Anonymous ACL submission

## Abstract

This paper introduces a new taxonomy of multilingual alignment for English-centric language models through token perturbation techniques. We propose two methods within this paradigm: the Language-Aware Token Boosting (LATB), which directly adds perturbations to desired language tokens, and its adaptive variant, the Adaptive Language-Aware Token Boosting (Adaptive-LATB), which dynamically adjusts perturbations based on the model’s confidence in the intended language. Extensive experiments show that our methods effectively enhance multilingual alignment by reducing language confusion and marginally improving summarization quality without requiring additional fine-tuning. Our code is publicly available<sup>1</sup>.

## 1 Introduction

Large Language Models (LLMs) have shown impressive performance, but their English-centric development limits their effectiveness for non-English users (Hadi et al., 2024, 2023). Recent efforts (Xue et al., 2021; Workshop et al., 2023; Wei et al., 2023) aim to enhance multilingual capabilities, though English-centric models still underperform in low-resource languages (Qin et al., 2024; OpenAI et al., 2024). One of the key issues is language confusion (Devine, 2024), where models fail to consistently generate the desired language, particularly in non-English contexts (Marchisio et al., 2024). Techniques to mitigate this include temperature lowering, few-shot prompting, and fine-tuning (Marchisio et al., 2024), but these come with limitations such as reduced responses diversity (Agarwal et al., 2024; Renze and Guven, 2024) or increased computational costs.

We propose a novel tuning-free paradigm for multilingual alignment, using perturbations di-

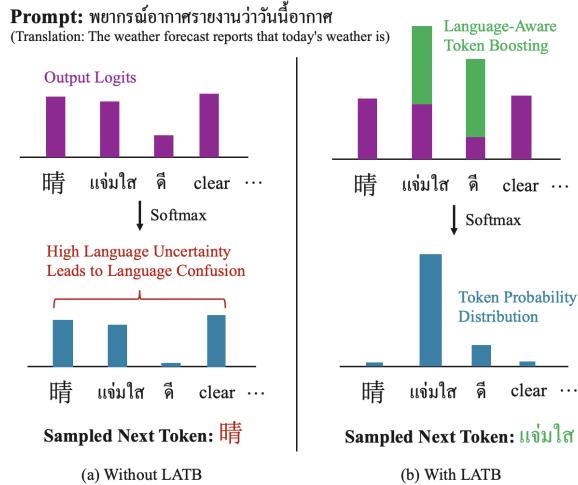


Figure 1: Language-Aware Token Boosting (LATB) enhances target language generation confidence by selectively boosting target language tokens.

rectly on the logits. This approach eliminates the need for fine-tuning and aligns the model’s outputs with the desired language, incurring minimal additional computational costs during inference. We introduce two methods within this paradigm: **Language-Aware Token Boosting (LATB)**, which applies language-specific token perturbations, and **Adaptive Language-Aware Token Boosting (Adaptive-LATB)**, which adapts perturbations by introducing perturbations selectively—only when the LLM exhibits uncertainty in generating one language over another.

We evaluate our methods on the XLSUM multilingual summarization benchmark (Hasan et al., 2021) across eight languages. Both LATB and Adaptive-LATB effectively reduce language confusion and enhance summarization performance compared to their respective base models and the multilingual-tuned model. We also analyze the effects of hyperparameters, including perturbation values and confidence difference thresholds.

In summary, our contributions are as follows:

<sup>1</sup><https://anonymous.4open.science/r/Language-Aware-Token-Boosting-Anonymous-7181>

- 062 1. We propose a novel tuning-free multilingual  
063 alignment paradigm based on logits pertur-  
064 bation, introducing two methods: LATB and  
065 Adaptive-LATB.
- 066 2. We evaluate our methods on the XLSUM  
067 benchmark, showing reduced language confu-  
068 sion and improved summarization quality.
- 069 3. We provide an analysis of the impact of hy-  
070 perparameters on the language confusion and  
071 summarization quality.

## 072 2 Related Work

**Multilingual Large Language Models.** Multilingual Large Language Models (MLLMs) are designed to process multiple languages simultaneously. The approaches for developing and optimizing these models can be broadly categorized into two main types: parameter-tuning alignment (PTA) and parameter-frozen alignment (PFA) (Qin et al., 2024). The PTA approach involves tuning the model’s parameters to enable multilingual capabilities. This tuning can occur at various stages, including pretraining (Xue et al., 2021; Chowdhery et al., 2022; Workshop et al., 2023; Jiang et al., 2023, 2024), supervised fine-tuning (SFT) (Chung et al., 2022; Muennighoff et al., 2023; Devine, 2024; Pipatanakul et al., 2023), reinforcement learning with human feedback (RLHF) (Lai et al., 2023b; Touvron et al., 2023; GLM et al., 2024; Bai et al., 2023), and downstream task fine-tuning (Lepikhin et al., 2020; Rosenbaum et al., 2022). In contrast, PFA methods do not require parameter tuning for multilingual alignment. Instead, they primarily rely on prompting techniques (Abdelali et al., 2024; Winata et al., 2023; Lu et al., 2024; Puduppully et al., 2023) and retrieval-augmented alignment (He et al., 2023; Zhang et al., 2023; Conia et al., 2023). Our proposed method falls within the PFA category. To the best of our knowledge, our study is the first to introduce a new taxonomy for logits perturbation-based multilingual alignment.

**Language Confusion.** Language confusion refers to the inconsistent ability of LLMs to generate responses in a target language. This phenomenon has been observed across various NLP tasks, such as machine translation (Vu et al., 2022; Li and Murray, 2023), summarization (Wang et al., 2023; Yu et al., 2022), and question answering (Holtermann et al., 2024). While this issue has been systematically studied with

various proposed methods mitigating it (Marchisio et al., 2024), our study introduces a novel and cost-effective approach to mitigate language confusion using token perturbation methods.

## 115 3 Approach

### 116 3.1 Token Language Identification

We identify tokens to boost based on the target language using a Unicode filtering method following (Wen-Yi and Mimno, 2023). Specifically, a token is considered valid if all its characters belong to the Unicode set defined for the target language. We also include numbers, special characters, and the end of sentence tokens in the desired set.

### 124 3.2 Perturbation Vector

We construct a perturbation vector,  $\mathbf{p}$ , based on the set of desired token indices  $I$ . Each element corresponding to an index in  $I$  is assigned a perturbation value  $\alpha \geq 0$ , as defined in Equation 1.

$$\mathbf{p}_i = \begin{cases} \alpha & \text{if } i \in I, \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

### 130 3.3 Logits Perturbation Methods

In this study, we explore two variants of the Logits Perturbation Method: LATB and Adaptive-LATB.

#### 133 3.3.1 Language-Aware Token Boosting 134 (LATB)

We introduce perturbations to the logits by adding a perturbation value  $\alpha$  to the selected logits to align them with the desired language. The method is detailed in Algorithm 1.

---

#### Algorithm 1 Vanilla LATB

---

```

logits ← LLM(x)
logits' ← logits + p      ▷ Logits Perturbation
y ← Softmax(logits')

```

---

#### 139 3.3.2 Adaptive Language-Aware Token 140 Boosting (Adaptive-LATB)

Adding logits in the vanilla LATB may suppress the ability to express tokens in another language when necessary. In contrast, the Adaptive LATB perturbs logits only when the LLM is not confident about the language it intends to express. The confidence difference threshold, controlled by the hyperparameter  $\beta$  ( $0 \leq \beta \leq 1$ ), determines the model’s confidence difference threshold in one language over another. This design enables the model

150 to switch languages when it is highly confident.  
151 The details of the Adaptive LATB algorithm are  
152 provided in Algorithm 2.

---

**Algorithm 2** Adaptive LATB

---

```
logits ← LLM(x)
y ← Softmax(logits)
a ← max({ $y_i \mid y_i \in \mathbf{y}$  and  $i \in I\})
b ← max({ $y_i \mid y_i \in \mathbf{y}$  and  $i \notin I\})
if  $|a - b| < \beta$  then
    logits' ← logits + p ▷ Logits Perturbation
    y ← Softmax(logits')
end if$$ 
```

---

## 4 Evaluation Metrics

We evaluate the model based on two key aspects: *Language Confusion*, which measures the model’s misalignment with the target language, and *Performance*, which assesses the quality of the generated summaries.

### 4.1 Language Confusion Metrics

We evaluate language confusion at three distinct levels to capture both fine-grained and overall effects: token-level, line-level, and response-level language confusion.

**Token-level Language Confusion.** We determine each token’s language based on its Unicode and calculate token-level misalignment rates for each response. These rates are then averaged across all responses to report the final metric.

**Line-level Language Confusion.** We segment each response by line and utilize an off-the-shelf language identification (LID) tool, FastText (Grave et al., 2018), to determine the language of each line. We calculate the average language misalignment per response and report the overall average across all responses.

**Response-level Language Confusion.** We input the entire response into the FastText (Grave et al., 2018) language identification and calculate the average language misalignment across all responses, reporting this as the final metric.

### 4.2 Performance Metrics

We assess summarization performance using three widely adopted metrics: ROUGE-1, ROUGE-2, and ROUGE-L (Lin, 2004). These metrics evaluate the overlap of unigrams, bigrams, and longest

common subsequences, respectively, between the generated summaries and the reference summaries.

## 5 Experiments

**Models.** We use Llama3 8B Instruct (Lai et al., 2023a) as the base English-centric model. To assess our method’s effectiveness, we compare it against Suzume 8B Multilingual (Devine, 2024), a multilingual fine-tuned version of Llama3 8B Instruct trained on a multilingual conversational dataset.

**Benchmark.** We adopt the multilingual summarization XLSUM dataset (Hasan et al., 2021) as the benchmark for our evaluation. This dataset is particularly suitable for our study as it allows models to generate extended responses, which can be systematically evaluated using quantitative metrics.

In our study, we select 4 High Resource Languages (HRL): Russian (ru), Simplified Chinese (zh), Japanese (ja), and French (fr), as well as 4 Medium Resource Languages (MRL): Korean (ko), Thai (th), Hindi (hi), and Arabic (ar). The categorization of languages follows (Lai et al., 2023a). For each language, we sample up to 1,000 examples for evaluation.

**Language Confusion Results.** We compare our methods with Llama3 8B Instruct (Grattafiori et al., 2024) using both standard and strict prompts that explicitly instruct the model to generate responses in a target language, and with Suzume 8B Multilingual (Devine, 2024). Our methods demonstrate effectiveness in reducing language confusion compared to its base model. Furthermore, our methods outperform the multilingual fine-tuned model, highlighting their effectiveness in reducing language confusion without incurring the cost of fine-tuning. The language confusion results are presented in Table 1.

**Summarization Quality Results.** The results reported in Table 2 demonstrate that our methods generate higher-quality responses compared to both the Llama3 baseline model (Grattafiori et al., 2024) and the Suzume 8B Multilingual model (Devine, 2024) without additional fine-tuning requirements.

## 6 Analysis

**Impact of Hyperparameters.** In LATB, increasing  $\alpha$  reduces language confusion, with ROUGE scores rising initially before declining. At the  $\alpha$  yielding the highest ROUGE scores, the model balances effectively expressing technical terms in

Table 1: Language confusion across different methods evaluated on eight languages, reported as Token-level/Line-level/Response-level language confusion in percentage.

	Llama3 8B-I	Llama3 8B-I (Strict Prompt)	Suzume 8B-Multilingual	Llama3 8B-I + LATB ( <i>Ours</i> )	Llama3 8B-I + Adaptive LATB ( <i>Ours</i> )
High Resource Languages (HRL)					
ru	80.66/93.83/92.50	5.02/4.10/2.90	3.04/2.30/2.10	<b>0.28/0.44/0.10</b>	0.48/ <b>0.38/0.10</b>
zh	85.92/98.90/98.90	14.17/9.69/9.10	7.56/0.89/0.90	<b>4.78/0.00/0.00</b>	5.37/0.10/0.00
ja	85.10/98.83/98.31	10.15/9.29/4.16	5.96/0.73/0.67	<b>3.51/0.70/0.11</b>	4.05/ <b>0.16/0.11</b>
fr	0.24/47.14/40.6	0.26/0.39/ <b>0.20</b>	0.31/0.37/0.30	<b>0.11/0.35/0.20</b>	0.18/ <b>0.25/0.30</b>
Medium Resource Languages (MRL)					
ko	85.46/99.60/99.63	16.72/30.79/27.27	8.28/11.74/12.36	<b>3.45/9.98/10.36</b>	4.56/10.12/11.45
th	86.03/99.80/99.39	3.67/9.80/2.30	2.16/1.16/0.84	0.43/0.18/ <b>0.00</b>	<b>0.38/0.00/0.00</b>
hi	86.18/99.05/98.50	1.67/8.59/0.40	2.77/3.36/2.50	<b>0.23/0.74/0.10</b>	0.26/0.89/ <b>0.00</b>
ar	86.21/99.27/98.30	9.98/11.94/5.60	5.63/2.95/2.60	<b>0.37/0.28/0.00</b>	0.54/ <b>0.22/0.00</b>

Table 2: Summarization performance across different methods evaluated on eight languages, reported as ROUGE-1/ROUGE-2/ROUGE-L in percentage.

	Llama3 8B-I	Llama3 8B-I (Strict Prompt)	Suzume 8B-Multilingual	Llama3 8B-I + LATB ( <i>Ours</i> )	Llama3 8B-I + Adaptive LATB ( <i>Ours</i> )
High Resource Languages (HRL)					
ru	4.89/0.96/4.18	20.44/9.26/13.41	19.35/8.32/12.42	20.83/ <b>9.46/13.60</b>	<b>21.00/9.42/13.58</b>
zh	0.80/0.32/0.69	19.41/8.99/13.73	19.31/8.59/13.38	<b>20.70/9.44/14.64</b>	20.55/9.28/14.52
ja	26.42/12.53/16.84	26.48/12.43/16.97	26.13/11.73/16.55	27.54/ <b>12.95/17.70</b>	<b>27.89/12.92/17.89</b>
fr	14.71/6.09/10.49	19.98/8.90/13.71	18.56/7.89/12.47	<b>20.13/9.05/13.74</b>	19.97/8.89/13.59
Medium Resource Languages (MRL)					
ko	2.27/0.24/2.12	14.66/6.14/10.16	15.30/6.13/10.47	16.41/6.78/11.38	<b>16.88/7.03/11.67</b>
th	1.79/0.45/1.51	29.24/13.99/15.62	28.99/13.29/15.14	29.77/14.07/15.79	<b>30.97/14.74/16.41</b>
hi	0.93/0.36/0.73	29.83/16.41/19.03	27.71/14.78/17.52	29.68/16.41/19.00	<b>29.77/16.41/19.05</b>
ar	1.54/0.19/1.42	19.22/7.46/11.66	19.60/7.09/11.69	<b>20.44/8.02/12.45</b>	19.79/7.62/11.84

English while minimizing language confusion at both line and response levels. Beyond this optimal point, higher  $\alpha$  values suppress tokens in non-target languages, leading to a performance drop.

For Adaptive LATB, higher  $\beta$  values also reduce language confusion, with ROUGE scores improving slightly until an inflection point. Excessively high  $\beta$  values hinder the model’s ability to generate tokens in non-target languages, resulting in a slight performance decline. The effect of hyperparameters are illustrated in Appendix C.

**Performance Improvements.** Our analysis highlights a strong correlation between performance improvements from LATB and the degree of language confusion without LATB. This finding suggests that language confusion contributes to performance degradation. By incorporating LATB, we effectively mitigate this issue, leading to performance gains. The relationship is illustrated in Figure 6 in Appendix D.

**Vanilla vs. Adaptive LATB.** Both methods deliver comparable performance. However, Vanilla LATB requires an optimal hyperparameter search to produce non-target language output when needed while accurately generating results in the target language. In contrast, Adaptive LATB is less sensitive to hyperparameters and supports non-

target language output as required.

## 7 Conclusion

This paper introduces a novel approach to multilingual alignment for English-centric language models through token perturbation techniques. We proposed the Language-Aware Token Boosting (LATB) and its adaptive variant, Adaptive-LATB. Extensive experiments demonstrate that our methods significantly reduce language confusion compared to base model and outperform its multilingual fine-tuned model. This highlights the efficiency and practicality of our approach for enhancing multilingual language model capabilities.

## Limitations and Future Work

Our work shows promising results but has several limitations. First, the methods struggle with aligning LLMs to untrained or out-of-vocabulary (OOV) tokens. Second, reliance on Unicode-based language identification is less effective for languages with significant overlap with Latin scripts. Finally, hyperparameter tuning is needed to balance language confusion and multilingual expression. Future work could improve OOV token handling, develop better token-based language identification techniques, and design language-agnostic hyperparameter selection methods.

## References

- Ahmed Abdelali, Hamdy Mubarak, Shammur Absar Chowdhury, Maram Hasanain, Basel Mousi, Sabri Boughorbel, Yassine El Kheir, Daniel Izham, Fahim Dalvi, Majd Hawasly, et al. 2024. [Larabench: Benchmarking arabic ai with large language models](#). *Preprint*, arXiv:2305.14982.
- Arav Agarwal, Karthik Mittal, Aidan Doyle, Pragnya Sridhar, Zipiao Wan, Jacob Arthur Doughty, Jaromir Savelka, and Majd Sakr. 2024. [Understanding the role of temperature in diverse question generation by gpt-4](#). In *Proceedings of the 55th ACM Technical Symposium on Computer Science Education V. 2*, SIGCSE 2024, page 1550–1551. ACM.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. [Qwen technical report](#). *Preprint*, arXiv:2309.16609.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. [Palm: Scaling language modeling with pathways](#). *Preprint*, arXiv:2204.02311.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. [Scaling instruction-finetuned language models](#). *Preprint*, arXiv:2210.11416.
- Simone Conia, Min Li, Daniel Lee, Umar Farooq Minhas, Ihab Ilyas, and Yunyao Li. 2023. [Increasing coverage and precision of textual information in multilingual knowledge graphs](#). *Preprint*, arXiv:2311.15781.
- Peter Devine. 2024. [Tagengo: A multilingual chat dataset](#). *Preprint*, arXiv:2405.12612.
- Team GLM, ;, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, et al. 2024. [Chatglm: A family of large language models from glm-130b to glm-4 all tools](#). *Preprint*, arXiv:2406.12793.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Edouard Grave, Piotr Bojanowski, Prakhar Gupta, Armand Joulin, and Tomas Mikolov. 2018. Learning word vectors for 157 languages. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)*.
- Muhammad Usman Hadi, Qasem Al Tashi, Abbas Shah, Rizwan Qureshi, Amgad Munee, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, et al. 2024. Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects. *Authorea Preprints*.
- Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili, et al. 2023. A survey on large language models: Applications, challenges, limitations, and practical usage. *Authorea Preprints*.
- Tahmid Hasan, Abhik Bhattacharjee, Md Saiful Islam, Kazi Samin, Yuan-Fang Li, Yong-Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. 2021. [XL-sum: Large-scale multilingual abstractive summarization for 44 languages](#). *Preprint*, arXiv:2106.13822.
- Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. 2023. [Exploring human-like translation strategy with large language models](#). *Preprint*, arXiv:2305.04118.
- Carolin Holtermann, Paul Röttger, Timm Dill, and Anne Lauscher. 2024. [Evaluating the elementary multilingual capabilities of large language models with multiq](#). *Preprint*, arXiv:2403.03814.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. [Mistral 7b](#). *Preprint*, arXiv:2310.06825.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. [Mixtral of experts](#). *Preprint*, arXiv:2401.04088.
- Viet Dac Lai, Nghia Trung Ngo, Amir Pouran Ben Veyeh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Huu Nguyen. 2023a. [Chatgpt beyond english: Towards a comprehensive evaluation of large language models in multilingual learning](#). *Preprint*, arXiv:2304.05613.
- Viet Dac Lai, Chien Van Nguyen, Nghia Trung Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan A. Rossi, and Thien Huu Nguyen. 2023b. [Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback](#). *Preprint*, arXiv:2307.16039.
- Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. 2020. [Gshard: Scaling giant models with conditional computation and automatic sharding](#). *Preprint*, arXiv:2006.16668.
- Tianjian Li and Kenton Murray. 2023. Why does zero-shot cross-lingual generation fail? an explanation and a solution. *Preprint*, arXiv:2305.17325.

395	Chin-Yew Lin.	2004.	ROUGE: A package for automatic evaluation of summaries.	In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain.	Association for Computational Linguistics.	450
396						451
397						452
398						453
399	Hongyuan Lu, Haoran Yang, Haoyang Huang, Dong-dong Zhang, Wai Lam, and Furu Wei.	2024.	Chain-of-dictionary prompting elicits translation in large language models.	<i>Preprint</i> , arXiv:2305.06575.		454
400						455
401						456
402						457
403	Kelly Marchisio, Wei-Yin Ko, Alexandre Bérard, Théo Dehaze, and Sebastian Ruder.	2024.	Understanding and mitigating language confusion in llms.	<i>Preprint</i> , arXiv:2406.20052.		458
404						459
405						460
406						461
407	Niklas Muenninghoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al.	2023.	Crosslingual generalization through multitask finetuning.	<i>Preprint</i> , arXiv:2211.01786.		462
408						463
409						464
410						465
411						466
412						467
413	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, et al.	2024.	Gpt-4 technical report.	<i>Preprint</i> , arXiv:2303.08774.		468
414						469
415						470
416						471
417						472
418	Kunat Pipattanakul, Phatrasek Jirabovonvisut, Potsawee Manakul, Sittipong Sripaisarnmongkol, Ruangsak Patomwong, Pathomporn Chokchainant, and Kasima Tharnpipitchai.	2023.	Typhoon: Thai large language models.	<i>Preprint</i> , arXiv:2312.13951.		473
419						474
420						475
421						476
422						477
423	Ratish Puduppully, Anoop Kunchukuttan, Raj Dabre, Ai Ti Aw, and Nancy F. Chen.	2023.	Decomposed prompting for machine translation between related languages using large language models.	<i>Preprint</i> , arXiv:2305.13085.		478
424						479
425						480
426						481
427						482
428	Libo Qin, Qiguang Chen, Yuhang Zhou, Zhi Chen, Yinghui Li, Lizi Liao, Min Li, Wanxiang Che, and Philip S. Yu.	2024.	Multilingual large language model: A survey of resources, taxonomy and frontiers.	<i>Preprint</i> , arXiv:2404.04925.		483
429						484
430						485
431						486
432						487
433	Matthew Renze and Erhan Guven.	2024.	The effect of sampling temperature on problem solving in large language models.	<i>Preprint</i> , arXiv:2402.05201.		488
434						489
435						490
436	Andy Rosenbaum, Saleh Soltan, Wael Hamza, Yannick Versley, and Markus Boese.	2022.	Linguist: Language model instruction tuning to generate annotated utterances for intent classification and slot tagging.	<i>Preprint</i> , arXiv:2209.09900.		491
437						492
438						493
439						494
440						495
441	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill,					496
442						497
443						498
444						499
445						500
446						501
447						502
448						503
449						504
450	Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Dowd, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Cogan, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kociský, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruba, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev.	2024.	Gemma 2: Improving open language models at a practical size.	<i>Preprint</i> , arXiv:2408.00118.		505
451						506
452						507
453						508
454						509
455						510
456						511
457	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al.	2023.	Llama 2: Open foundation and fine-tuned chat models.	<i>Preprint</i> , arXiv:2307.09288.		512

512 Tu Vu, Aditya Barua, Brian Lester, Daniel Cer, Mo-  
513 hit Iyyer, and Noah Constant. 2022. Overcoming  
514 catastrophic forgetting in zero-shot cross-lingual gen-  
515 eration. *Preprint*, arXiv:2205.12647.

516 Jiaan Wang, Fandong Meng, Yunlong Liang, Tingyi  
517 Zhang, Jiarong Xu, Zhixu Li, and Jie Zhou. 2023.  
518 Understanding translationese in cross-lingual sum-  
519 marization. *Preprint*, arXiv:2212.07220.

520 Xiangpeng Wei, Haoran Wei, Huan Lin, Tianhao Li,  
521 Pei Zhang, Xingzhang Ren, Mei Li, Yu Wan, Zhi-  
522 wei Cao, Binbin Xie, et al. 2023. Polym: An open  
523 source polyglot large language model. *Preprint*,  
524 arXiv:2307.06018.

525 Andrea W Wen-Yi and David Mimno. 2023. Hyperpoly-  
526 glot LLMs: Cross-lingual interpretability in token  
527 embeddings. In *Proceedings of the 2023 Conference*  
528 *on Empirical Methods in Natural Language Process-*  
529 *ing*, pages 1124–1131, Singapore. Association for  
530 Computational Linguistics.

531 Genta Indra Winata, Alham Fikri Aji, Zheng-Xin  
532 Yong, and Thamar Solorio. 2023. The decades  
533 progress on code-switching research in nlp: A sys-  
534 tematic survey on trends and challenges. *Preprint*,  
535 arXiv:2212.09660.

536 BigScience Workshop, :, Teven Le Scao, Angela Fan,  
537 Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel  
538 Hesslow, Roman Castagné, Alexandra Sasha Luc-  
539 cioni, François Yvon, et al. 2023. Bloom: A 176b-  
540 parameter open-access multilingual language model.  
541 *Preprint*, arXiv:2211.05100.

542 Linting Xue, Noah Constant, Adam Roberts, Mihir  
543 Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua,  
544 and Colin Raffel. 2021. mt5: A massively multilingual  
545 pre-trained text-to-text transformer. *Preprint*,  
546 arXiv:2010.11934.

547 Sicheng Yu, Qianru Sun, Hao Zhang, and Jing Jiang.  
548 2022. Translate-train embracing translationese arti-  
549 facts. In *Proceedings of the 60th Annual Meeting of*  
550 *the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 362–370, Dublin, Ire-  
551 land. Association for Computational Linguistics.

553 Min Zhang, Limin Liu, Zhao Yanqing, Xiaosong Qiao,  
554 Su Chang, Xiaofeng Zhao, Junhao Zhu, Ming Zhu,  
555 Song Peng, Yinglu Li, et al. 2023. Leveraging mul-  
556 tilingual knowledge graph to boost domain-specific  
557 entity translation of ChatGPT. In *Proceedings of Ma-*  
558 *chine Translation Summit XIX, Vol. 2: Users Track*,  
559 pages 77–87, Macau SAR, China. Asia-Pacific Asso-  
560 ciation for Machine Translation.

## 561 A Experiment Details

562 We generate responses using the Llama3 8B In-  
563 struct model (Grattafiori et al., 2024) on eight dif-  
564 ferent languages from the XLSUM dataset (Hasan  
565 et al., 2021). The prompts utilized for this exper-  
566 iment are detailed in Appendix B. All responses

are generated with the sampling parameters set to  
a temperature of 1.0 and a top- $p$  value of 1.0. For  
LATB, the perturbation value  $\alpha$  is set to 5. For  
Adaptive LATB, the perturbation value is set to  
 $\alpha = 1000$ , and the confidence difference threshold  
is set to  $\beta = 0.8$ .

## 573 B Prompt Templates

Language	Strict Prompt
ru	Пожалуйста, кратко изложите текст на русском языке. Текст: {} Резюме:
zh	请用中文（简体）总结文本。文本: {} 总结:
ja	テキストを日本語で要約してください。テキスト: {} 要約:
fr	Veuillez résumer le texte en français. Texte : {} Résumé :
ko	텍스트를 한국어로 요약해 주세요. 텍스트: {} 요약:
th	กรุณารวบเรียงภาษาไทย ข้อความ: {} สรุป:
hi	कृपया पाठ का सारांश हिंदी में दों पाठ: {} सारांश:
ar	يرجى تلخيص النص باللغة العربية. النص: {} الملخص:

Figure 2: Strict prompt templates used in the experiment

Language	Standard Prompt
ru	Пожалуйста, кратко изложите текст. Текст: {} Резюме:
zh	请总结文本。文本: {} 总结
ja	テキストを要約してください。テキスト: {} 要約:
fr	Veuillez résumer le texte. Texte : {} Résumé :
ko	텍스트를 요약해 주세요. 텍스트: {} 요약:
th	กรุณารวบเรียงภาษาไทย ข้อความ: {} สรุป:
hi	कृपया पाठ का सारांश दों पाठ: {} सारांश:
ar	يرجى تلخيص النص. النص: {} الملخص:

Figure 3: Standard prompt templates used in the experiment

We design language-specific prompt templates to  
ensure consistency and adaptability across different  
languages during text generation. Each template  
provides a structured format where {} is replaced  
by the input text to summarize. The strict prompt  
templates include instructions to ensure the model  
generates output in the target language, whereas  
the standard prompt templates do not. The standard  
and strict prompt templates are shown in Figures 3  
and 2, respectively.

## 584 C Impacts of Hyperparameters

In LATB, the perturbation parameter  $\alpha$  influences  
the generated responses. To analyze its impact, we  
varied  $\alpha$  from 0 to 50 and recorded the correspond-  
ing results. These results are presented in Figure 4,  
illustrating the effect of  $\alpha$  on language confusion  
and summarization quality.

591 In Adaptive-LATB, we investigated the influence  
592 of the confidence difference threshold, denoted as  
593  $\beta$ . Specifically, we varied  $\beta$  from 0 to 0.9 while  
594 keeping the perturbation value  $\alpha$  fixed at 1000. The  
595 outcomes of this experiment are visualized in Fig-  
596 ure 5, highlighting how changes in  $\beta$  affect lan-  
597 guage confusion and summarization quality.

598 All responses across the experiments were gener-  
599 ated using a temperature setting of 1.0 and a top- $p$   
600 value of 1.0, ensuring consistency in sampling pa-  
601 rameters throughout the evaluations.

## 602 D Performance Improvements

603 The relationship is illustrated in Figure 6, which  
604 demonstrates a strong correlation between perfor-  
605 mance improvements using LATB and language  
606 confusion.

## 607 E Results on Multilingual Language 608 Model

609 We extend our experiments to the multilingual lan-  
610 guage model Gemma 2B Instruct (Team et al.,  
611 2024). Our methods reduce language confusion  
612 and slightly enhance summarization performance  
613 in most tested languages, as shown in Tables 3  
614 and 4, respectively.

## 615 F Vanilla vs. Adaptive LATB Example

616 As discussed, Vanilla LATB is sensitive to hyperpa-  
617 rameters, which can lead to a constraint on single-  
618 target language generation when  $\alpha$  is too high. In  
619 contrast, Adaptive LATB is less sensitive to hy-  
620 perparameters. The example output is shown in  
621 Figure 7. When  $\alpha$  is too high, Vanilla LATB gen-  
622 erates a Thai-dubbed version of English, whereas  
623 Adaptive LATB uses English directly, resulting in  
624 a more natural output.

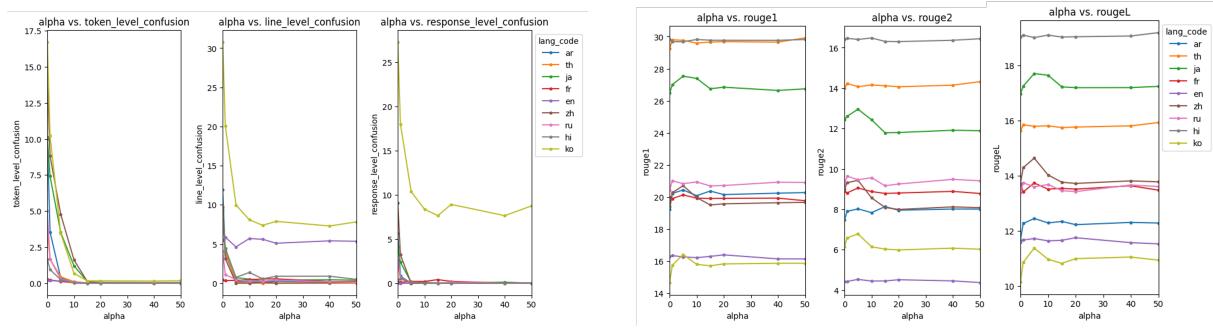


Figure 4: Impact of the Perturbation Value  $\alpha$  on Language Confusion and Performance in LATB

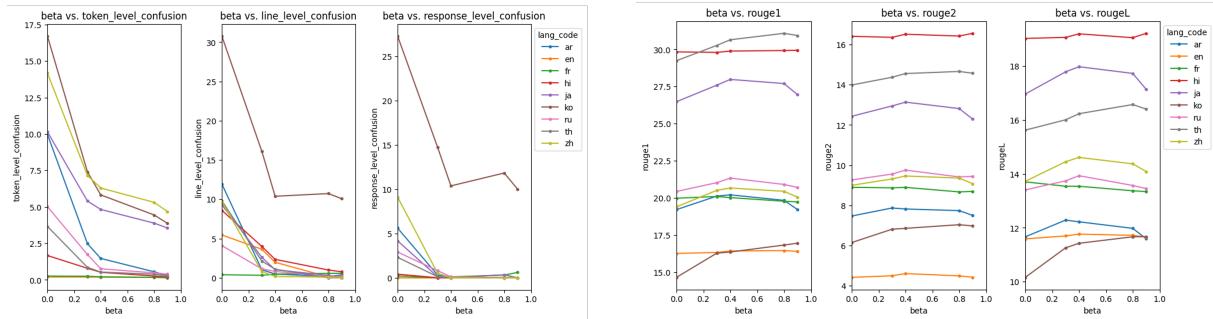


Figure 5: Impact of the Confidence Difference Threshold  $\beta$  on Language Confusion and Performance in Adaptive-LATB with  $\alpha$  fixed at 1000

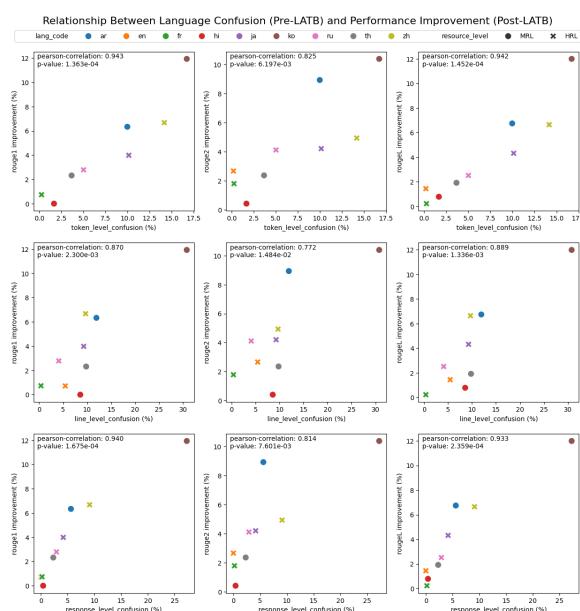


Figure 6: Performance improvements with LATB correlate strongly with language confusion levels

Figure 7: Output examples of Vanilla LATB with excessively high  $\alpha$  and Adaptive LATB

Table 3: Language confusion across different methods on Gemma model evaluated on eight languages, reported as Token-level/Line-level/Response-level language confusion in percentage.

	Gemma 2B-I	Gemma 2B-I (Strict Prompt)	Gemma 2B-I + LATB ( <i>Ours</i> )	Gemma 2B-I + Adaptive LATB ( <i>Ours</i> )
High Resource Languages (HRL)				
ru	4.65/2.94/2.90	3.78/2.30/2.30	<b>0.55/0.15/0.10</b>	0.66/0.20/0.20
zh	7.47/0.15/0.10	7.29/ <b>0.05/0.00</b>	<b>6.53/0.10/0.10</b>	6.62/0.10/0.10
ja	5.48/0.17/0.11	<b>5.53/0.00/0.00</b>	4.70/ <b>0.00/0.00</b>	<b>4.39/0.00/0.00</b>
fr	0.16/0.43/0.20	0.16/ <b>0.03/0.00</b>	0.89/0.05/ <b>0.00</b>	<b>0.12/0.15/0.00</b>
Medium Resource Languages (MRL)				
ko	6.99/11.86/11.82	6.90/9.87/10.00	<b>4.70/0.00/0.00</b>	5.87/7.27/7.09
th	1.47/0.71/0.48	1.50/0.57/0.36	<b>0.51/0.00/0.00</b>	<b>0.50/0.00/0.00</b>
hi	2.62/2.20/2.20	1.09/1.00/0.70	<b>0.24/0.00/0.00</b>	0.27/ <b>0.00/0.00</b>
ar	6.41/4.00/4.00	5.24/2.64/2.60	<b>0.18/0.03/0.00</b>	0.35/ <b>0.00/0.00</b>

Table 4: Summarization performance on Gemma model across different methods evaluated on eight languages, reported as ROUGE-1/ROUGE-2/ROUGE-L in percentage.

	Gemma 2B-I	Gemma 2B-I (Strict Prompt)	Gemma 2B-I + LATB ( <i>Ours</i> )	Gemma 2B-I + Adaptive LATB ( <i>Ours</i> )
High Resource Languages (HRL)				
ru	21.59/8.29/14.82	21.54/8.22/14.84	21.86/8.49/ <b>15.16</b>	<b>21.93/8.54/15.15</b>
zh	25.90/11.34/18.99	25.91/11.48/19.00	26.01/11.54/19.21	<b>26.14/11.60/19.29</b>
ja	30.28/12.68/20.08	30.08/12.64/20.16	<b>30.66/13.01/20.60</b>	29.86/12.53/20.26
fr	24.20/9.51/17.18	<b>24.39/9.61/17.32</b>	24.11/9.49/17.26	23.98/9.53/17.21
Medium Resource Languages (MRL)				
ko	<b>22.27/8.98/16.35</b>	22.11/8.81/15.88	21.84/8.87/15.90	21.76/8.65/16.07
th	31.59/13.99/17.19	30.97/13.55/17.01	31.44/13.75/17.15	<b>31.60/14.02/17.31</b>
hi	36.85/17.39/21.75	<b>37.19/17.52/22.08</b>	36.02/17.20/21.36	37.08/ <b>17.61/22.12</b>
ar	22.26/7.69/14.39	22.54/7.76/14.58	<b>22.66/8.03/14.62</b>	22.64/7.81/14.58