

# Cross-Lingual Pitfalls: Automatic Probing Cross-Lingual Weakness of Multilingual Large Language Models

Zixiang Xu<sup>1\*</sup>, Yanbo Wang<sup>1\*</sup>, Yue Huang<sup>2\*</sup>,  
Xiuying Chen<sup>1†</sup>, Jieyu Zhao<sup>3</sup>, Meng Jiang<sup>2</sup>, Xiangliang Zhang<sup>2†</sup>

<sup>1</sup>MBZUAI, <sup>2</sup>University of Notre Dame, <sup>3</sup>University of Southern California

Correspondence: xiuying.chen@mbzuai.ac.ae, xzhang33@nd.edu

## Abstract

Large Language Models (LLMs) have achieved remarkable success in Natural Language Processing (NLP), yet their cross-lingual performance consistency remains a significant challenge. This paper introduces a novel methodology for efficiently identifying inherent cross-lingual weaknesses in LLMs. Our approach leverages beam search and LLM-based simulation to generate bilingual question pairs that expose performance discrepancies between English and target languages. We construct a new dataset of over 6,000 bilingual pairs across 16 languages using this methodology, demonstrating its effectiveness in revealing weaknesses even in state-of-the-art models. The extensive experiments demonstrate that our method precisely and cost-effectively pinpoints cross-lingual weaknesses, consistently revealing over 50% accuracy drops in target languages across a wide range of models. Moreover, further experiments investigate the relationship between linguistic similarity and cross-lingual weaknesses, revealing that linguistically related languages share similar performance patterns and benefit from targeted post-training. Code is available at <https://github.com/xzx34/Cross-Lingual-Pitfalls>.

## 1 Introduction

Large language models (LLMs) have rapidly ascended to prominence in Natural Language Processing (NLP), gaining recognition for their exceptional performance across various tasks, spanning from the sciences (Li et al., 2024a; Guo et al., 2023; Huang et al., 2024c; Wang et al., 2025b; Xu et al., 2025a) to the development of LLM-based agents (Liu et al., 2023b, 2024b). Recent advancements have driven research on enhancing LLMs' multilingual capabilities (Zhao et al., 2024a; Wang et al.,

\*Equal contribution.

†Corresponding authors.

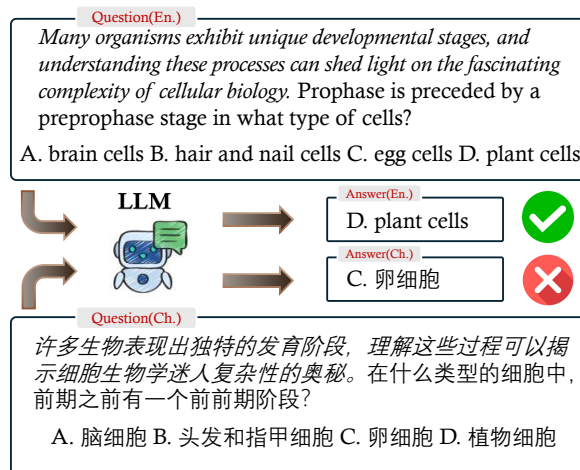


Figure 1: An example of an English-Chinese question pair discovered by our search methodology (where the Chinese question is semantically equivalent to the English) highlights the cross-lingual performance gap: even GPT-4o, despite its strong multilingual capabilities, provides the correct answer in English but gives an incorrect response in Chinese.

2025d), improving their effectiveness in addressing real-world problems with greater nuance and global reach.

Despite advancements, inconsistencies in model performance across languages remain a significant challenge (Xu et al., 2024). The proficiency demonstrated in English often fails to generalize to other languages, resulting in errors in other linguistic contexts, as exemplified in Figure 1. To effectively enhance the cross-lingual consistency of these models, an initial and crucial step is the identification of their inherent cross-lingual weaknesses. Since English is the primary training language for LLMs, and they generally perform best in English (Li et al., 2024b), we define cross-lingual weakness in this paper as: *For a given question presented in multiple languages, a model answers correctly in English but incorrectly in at least one other language.* This definition requires the model to provide the correct answer in English, as failure across all languages

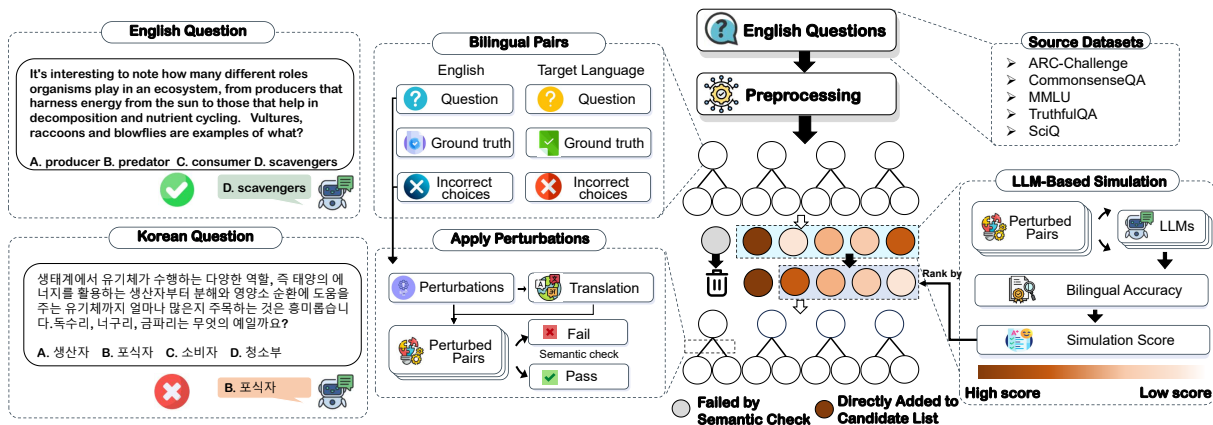


Figure 2: The overview of the proposed methodology for generating questions that precisely challenge the cross-lingual capabilities of LLMs. As depicted, the pipeline initiates with sampling English questions and creating bilingual pairs. Iterative perturbation, driven by a beam search strategy and guided by LLM-based simulation scores, refines these pairs to maximize performance divergence between English and the target language. The resulting candidate list of question pairs is designed to highlight inherent cross-lingual weaknesses in LLMs.

would likely indicate a knowledge-related limitation rather than a cross-lingual weakness.

To efficiently uncover these weaknesses, we propose a beam search-based methodology. This approach leverages existing, high-quality English datasets and iteratively introduces perturbations to the English questions. These perturbations are designed to increase question complexity and cognitive demand for comprehension and completion. The goal is to prevent models from generating answers based on superficial cues without genuine language understanding (Stacey et al., 2020; Bhargava et al., 2021), thereby exposing disparities in cross-lingual capabilities. Our approach begins with sourcing English questions from high-quality existing datasets, which are then translated into the target language to form bilingual question pairs. Crucially, this translation process incorporates a semantic check to guarantee that the meaning of the questions is preserved and the correct answer remains consistent across languages. Subsequently, these pairs undergo iterative perturbations, generating a diverse set of perturbed pairs. Then we employ an LLM-based simulation framework that assigns a simulation score measuring the effectiveness of revealing cross-lingual weaknesses, to each pair for ranking. The top-ranked pairs are iteratively perturbed to further refine the search process. Finally, question pairs with consistently high accuracy in English but significant performance drops in the target language are added to the candidate list to expose cross-lingual weaknesses in LLMs.

Furthermore, to study how cross-lingual weaknesses are relevant to linguistic similarity, we con-

ducted exploratory experiments. Our key findings reveal that: 1) languages closer in linguistic terms tend to share similar weaknesses; and 2) fine-tuning LLMs on one language improves performance more significantly in linguistically similar languages. These results highlight that linguistic relationships strongly influence cross-lingual performance.

In summary, our contributions are: 1) We present an efficient, precise methodology for identifying LLM cross-lingual weaknesses. 2) Based on the proposed methodology, we construct a novel, 16-language dataset with over 6,000 bilingual pairs to challenge cross-lingual capabilities. 3) Extensive experiments on the dataset quantitatively analyze the relationship between cross-lingual weaknesses and linguistic similarities and fine-tuning experiments demonstrate the potential for targeted cross-lingual improvement.

## 2 Methodology

In this section, we introduce our methodology for automatically probing the cross-lingual weakness of multilingual LLMs. As illustrated in Figure 2, our goal is to generate questions that precisely challenge the cross-lingual capabilities of LLMs by identifying cases where the model performs well in English but struggles with the same questions when presented in a specific target language.

### 2.1 Method Overview

To achieve the goal described above, we first sample a set of English questions and translate them into bilingual pairs, where each pair consists of an

English question and its counterpart in the target language. We then iteratively introduce perturbations to these pairs using a *beam search strategy*, guided by maximizing the accuracy discrepancy between English and the target language (i.e., to retain the accuracy on English questions but make accuracy on target language questions drop as much as possible). This search-and-perturbation approach is inspired by prior work on uncovering model vulnerabilities through optimization-guided example construction (Huang et al., 2025b). During beam search, a *LLM-based simulation* is utilized to guide the search process in identifying the model’s weaknesses in the target language. Based on the *search optimization strategies*, we aim to balance the trade-off between problem generation and computational cost. It ultimately produces a candidate list of English–target language question pairs, effectively highlighting the model’s cross-lingual weaknesses.

## 2.2 Problem Formulation

Let  $\mathcal{B} = \{(q_i^E, q_i^T)\}_{i=1}^W$  denote our original set of  $W$  bilingual pairs. Each bilingual pair  $(q^E, q^T)$  is formally represented as:

$$(q^E, q^T, a_\star^E, a_\star^T, \mathcal{A}_\neg^E, \mathcal{A}_\neg^T), \quad (1)$$

where  $q^E, q^T \in \mathcal{Q}$  represent question texts in English and the target language, respectively.  $a_\star^E, a_\star^T \in \mathcal{A}$  are the corresponding ground-truth answers.  $\mathcal{A}_\neg^E, \mathcal{A}_\neg^T$  denote incorrect answer choices.

During the beam search process, we iteratively apply perturbations to bilingual pairs. Specifically, given an English question  $q^E$  from a bilingual pair and an incorrect answer  $\alpha^E \in \mathcal{A}_\neg^E$ , the perturbation function  $\varphi$  generates a semantically irrelevant yet contextually plausible perturbation:

$$\delta q^E = \varphi(q^E, \alpha^E), \quad (2)$$

where  $\varphi : \mathcal{Q} \times \mathcal{A} \rightarrow \mathcal{Q}$  modifies  $q^E$  while preserving its original semantics and embedding patterns influenced by the incorrect answer  $\alpha^E$ . Here,  $\varphi$  is a proxy LLM utilized for adding the perturbation.

The perturbed English question is then formed as:

$$q^{E'} = \oplus(q^E, \delta q^E), \quad (3)$$

where  $\oplus : \mathcal{Q} \times \mathcal{Q} \rightarrow \mathcal{Q}$  denotes a context-sensitive insertion of the perturbation into  $q^E$ . During implementation,  $\oplus$  is a concatenation operation.

To maintain consistency across languages, the corresponding perturbation in the target language

is generated as  $\delta q^T = \mathcal{T}(\delta q^E)$ , where  $\mathcal{T} : \mathcal{Q} \rightarrow \mathcal{Q}$  is a translation module that strictly translates the inserted perturbation without modifying other parts of the question. This results in the perturbed target-language question:  $q^{T'} = \oplus(q^T, \delta q^T)$ .

We optimize the perturbation to minimize the model’s accuracy in the target language while maintaining near-perfect performance in English. Formally, our objective is:

$$\begin{aligned} \min_{\delta q^E} \quad & \mathbb{E} \left[ \mathbb{I}(\mathcal{F}(q^{T'}) = a_\star^T) \right] \\ \text{s.t.} \quad & \mathbb{E} \left[ \mathbb{I}(\mathcal{F}(q^{E'}) = a_\star^E) \right] \geq 1 - \epsilon \\ & \mathbb{S}(q^E, q^{E'}) \geq \theta, \quad \mathbb{S}(q^{E'}, q^{T'}) \geq \theta'. \end{aligned} \quad (4)$$

where  $\mathcal{F}$  represents the LLM’s response function,  $\mathbb{S}$  is a semantic similarity function,  $\theta$  and  $\theta'$  are threshold values ensuring semantic consistency, and  $\mathbb{I}(\cdot)$  is the indicator function, which returns 1 if the predicted answer is correct and 0 otherwise.

The first constraint ensures that perturbations  $\delta q^E$  preserve the model’s accuracy in English ( $\mathbb{E} \geq 1 - \epsilon$ ), while the second set of constraints ensures that the perturbed and original questions remain semantically equivalent in both English and the target language.

## 2.3 LLM-Based Simulation

LLM-based simulation utilizes a set of LLMs to answer perturbed questions and derive a simulation score based on the accuracy relationship between bilingual pairs.

The simulation employs a collection of LLMs, denoted as  $\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_K\}$ , to quantify the cross-lingual performance gap introduced by perturbations. For each perturbed bilingual pair  $(q^{E'}, q^{T'})$ , each  $\mathcal{M}_k \in \mathcal{M}$  generates predicted answers:

$$\hat{a}_k^{E'} = \mathcal{M}_k(q^{E'}), \quad \hat{a}_k^{T'} = \mathcal{M}_k(q^{T'}). \quad (5)$$

The correctness of these predictions is assessed by comparing them to the ground truth answers:

$$\beta_k^{E'} = \mathbb{I}(\hat{a}_k^{E'} = a_\star^E), \quad \beta_k^{T'} = \mathbb{I}(\hat{a}_k^{T'} = a_\star^T). \quad (6)$$

The average accuracy across all models is computed as:

$$\bar{\beta}^{E'} = \frac{1}{K} \sum_{k=1}^K \beta_k^{E'}, \quad \bar{\beta}^{T'} = \frac{1}{K} \sum_{k=1}^K \beta_k^{T'}. \quad (7)$$

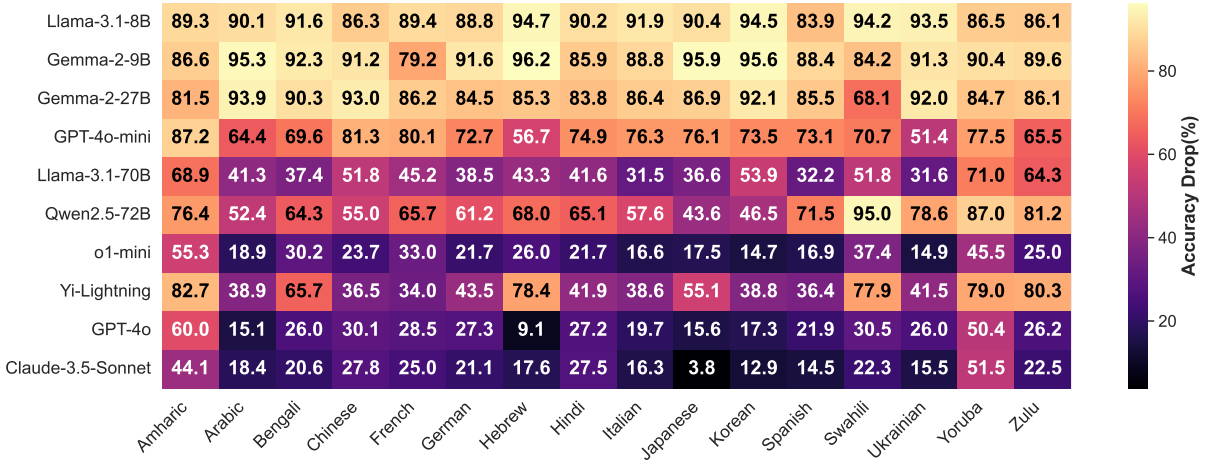


Figure 3: Evaluation of 10 models on our generated 6,600 bilingual pairs across 16 languages. While all models achieve nearly 100% accuracy in English, most experience an average accuracy drop of over 50% in the target languages. Even state-of-the-art multilingual models like GPT-4o and Claude-3.5-sonnet exhibit significant cross-lingual weaknesses.

To evaluate the effectiveness of each perturbation, we define a simulation score  $V(q^{E'}, q^{T'})$  that highlights significant performance discrepancies:

$$V(q^{E'}, q^{T'}) = \left(\bar{\beta}^{E'}\right)^\gamma - \bar{\beta}^{T'}, \quad (8)$$

where  $\gamma > 1$  is an exponent that amplifies high English accuracy. This formulation prioritizes bilingual pairs where the model maintains strong performance in English ( $\bar{\beta}^{E'} \approx 1$ ) but exhibits significant degradation in the target language ( $\bar{\beta}^{T'}$ ).

## 2.4 Beam Search with Optimization Strategies

Since beam search is an effective heuristic for exploring a constrained search space, we employ it to solve the objective function in Equation 4 by greedily identifying the top perturbation candidates produced by the proxy LLMs. These candidates are evaluated and ranked based on their effectiveness in causing performance discrepancies, defined as  $V(q^{E'}, q^{T'})$  in Equation 8.

Here, the search width  $w$  determines the number of top-ranked bilingual pairs retained after ranking at each iteration, effectively controlling the breadth of the search at each level of the search tree. The initial search depth  $d_1$  specifies the maximum depth of the search tree explored in the initial phase, corresponding to the maximum number of perturbation iterations applied to a question.

Next, we discuss key factors in the Beam Search process that determine: 1) when a perturbed question qualifies as a valid candidate, 2) when to terminate the search for a given bilingual pair, and 3) how to ensure diversity within the set of candidates.

**Inclusion Threshold Strategy.** A bilingual pair is immediately included in the candidate list if its simulation score exceeds a predefined inclusion threshold  $\theta_{inc}$ , ensuring early termination for critical perturbations. Otherwise, the top  $w$  scoring pairs are selected to advance to the next search level, where the search depth increments by one.

**Early Stopping Mechanism.** To adaptively adjust the search depth based on the quality of discovered perturbations, we introduce a potential threshold  $\theta_{pot}$ , which determines whether the search should continue beyond the initial depth. Specifically, the search depth  $d$  at iteration  $t$  is updated as follows:

$$d_t = \begin{cases} d_2, & \text{if } \max_{q^{E'}, q^{T'} \in \mathcal{B}_t} V(q^{E'}, q^{T'}) \geq \theta_{pot}, \\ d_1, & \text{otherwise.} \end{cases} \quad (9)$$

where  $\mathcal{B}_t$  represents the set of bilingual pairs at iteration  $t$ . The search process continues until reaching the maximum allowable depth,  $d_{max} = \max(d_1, d_2)$ . Thus, if at any iteration a perturbation achieves a simulation score surpassing  $\theta_{pot}$ , the search depth is expanded to  $d_2$ , allowing further exploration. Otherwise, the search remains at  $d_1$ . The process terminates when  $d_{max}$  is reached.

**Redundancy Control Mechanism.** To ensure diversity in the candidate list, if  $r$  bilingual pairs originating from the same initial question have already been included in the candidate list, all remaining bilingual pairs derived from that question are discarded from further exploration. This prevents excessive redundancy and ensures a wider variety of perturbed questions in the candidate list.

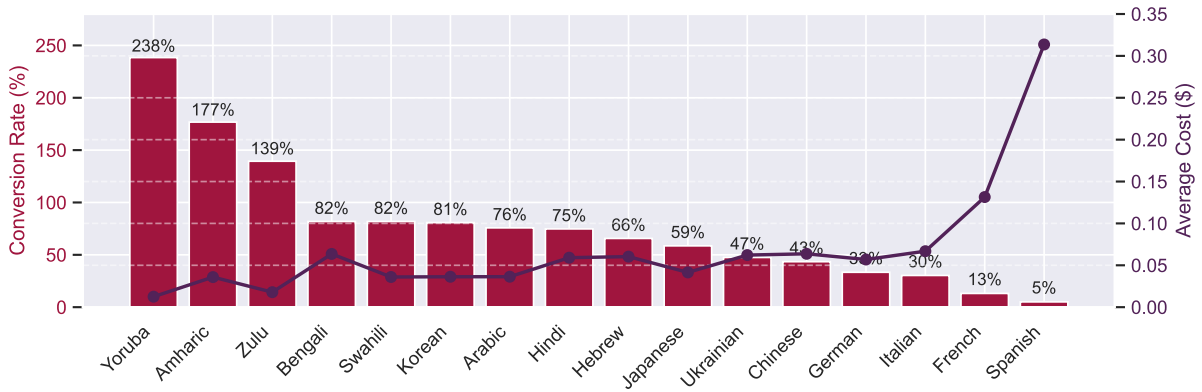


Figure 4: Analysis of question conversion rates and generation costs across 16 languages based on all pairs in our candidate list. The bar chart (red) shows question conversion rates for different languages, while the line chart (purple) represents cost of generating a single question. Notably, in most languages, identifying a bilingual pair that exposes cross-lingual weaknesses costs less than \$0.05. However, for languages structurally and lexically closer to English, such as French and Spanish, finding weaknesses becomes significantly harder, leading to higher costs.

### 3 Experiment

#### 3.1 Experiment Overview

In this section, we conduct a series of experiments to evaluate the effectiveness of our proposed method as well as to explore the cross-lingual weaknesses of multilingual models. Overall, we mainly aim to address the following questions:

- **RQ1:** How effective and efficient is our method in identifying the cross-lingual weaknesses of multilingual models? (§3.2)
- **RQ2:** Are the identified weaknesses language-specific? How can we understand the linguistic connection between cross-lingual weaknesses and the languages involved? (§3.3)
- **RQ3:** Furthermore, to what extent does language-specific fine-tuning enhance cross-lingual performance, and how is the fine-tuning improvement associate with the relationships of different languages? (§3.4)

#### 3.2 Cross-Lingual Weakness Identification

To answer RQ1, based on the proposed method, we generated initial bilingual pairs using GPT-4o and employed GPT-4o-mini for perturbation generation. Perturbations were translated using the Google Translate API (Google). We then employed  $W$  cost-effective multilingual models in  $\mathcal{M}$  for LLM-based simulation to generate a set of candidates over 6,000 question pairs spanning 16 languages. These pairs are then used to evaluate the performance of 10 different models. Detailed experimental settings and parameter configurations are provided in Appendix B.

**Our method effectively identifies cross-lingual**

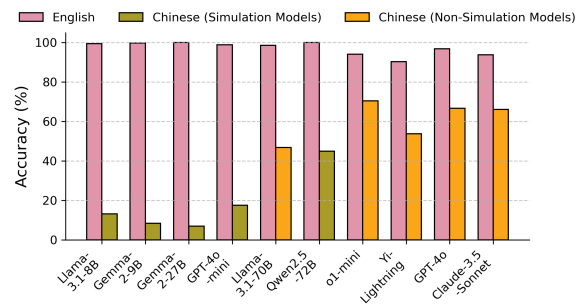


Figure 5: Performance of LLMs on our generated English-Chinese pairs. Even smaller models like Gemma-2-9B and Llama-3.1-8B achieve perfect accuracy in English, while more than half of the models score below 50% in Chinese. Despite their strong multilingual capabilities, GPT-4o and Claude-3.5-sonnet still exhibit over a 30% accuracy drop compared to English.

**weaknesses even in state-of-the-art models.** Taking Chinese as an example, we evaluated all models on our generated English-Chinese pairs and found that their accuracy dropped by nearly 60% on average when switching from English to Chinese, as shown in Figure 5. Notably, even the smallest models achieved perfect accuracy on English tasks (i.e., they have mastered the most knowledge of answering the questions), whereas the most advanced model, GPT-4o, still exhibited a substantial accuracy drop of nearly 30% in Chinese. Similar performance gaps were observed across other languages, as presented in Appendix B.2.

As shown in Figure 3, the accuracy drops across 16 languages highlight the cross-lingual performance gaps. Even Claude-3.5-sonnet experienced over 20% accuracy loss in most languages. This starkly demonstrates the effectiveness of our

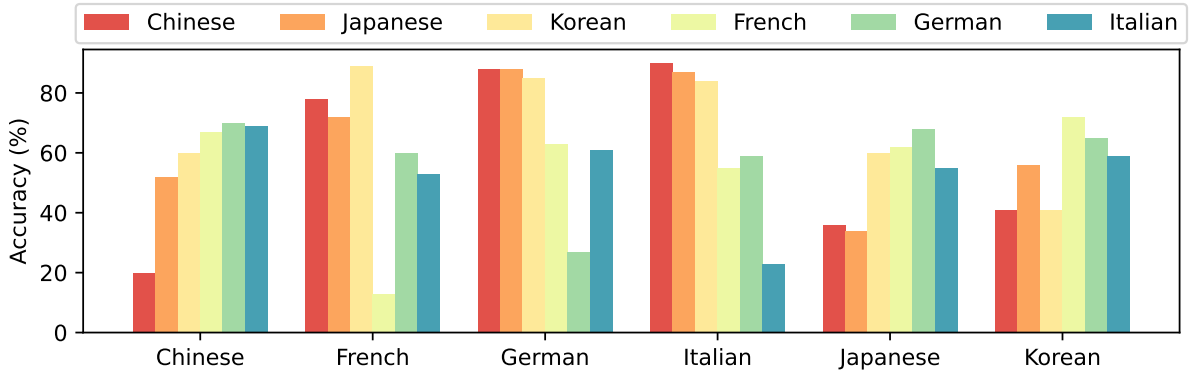


Figure 6: Accuracy of GPT-4o-mini on expanded bilingual pairs (Asian and European language families). The red bar represents accuracy for pairs expanded from Chinese seed pairs, while other colors show results for pairs expanded from other seed language pairs within these families.

method in identifying cross-lingual weaknesses in even state-of-the-art multilingual models.

Moreover, from Figure 5, we can observe that the models used for simulation typically exhibit greater accuracy degradation. By varying the models in  $\mathcal{M}$  for LLM-based simulation, we can discover specific cross-lingual weaknesses in any given LLMs. To investigate this, we replaced Gemma-2-27B and Qwen2.5-72B with GPT-4o in our simulation framework. A comparison between Figure 5 and Figure 10 reveals that: Qwen2.5-72B and Gemma-2-27B show minor accuracy improvements after being removed from the simulation models, GPT-4o—despite being a top-tier multilingual model—suffers a sharp 58% accuracy drop.

**Our method enables the cost-effective identification of cross-lingual weaknesses.** We evaluated the cost of generating bilingual pairs and analyzed the conversion rate for each language—i.e., the proportion of bilingual pairs successfully generated from an original English question—as illustrated in Figure 4. For most languages, the average cost of identifying a question that exposes a model’s cross-lingual weaknesses is as low as \$0.05.

Interestingly, for most languages, the cost of generating pairs is significantly lower, compared to the specific languages like French, Spanish, Italian, and German. This discrepancy can be explained by the greater linguistic similarities of these languages to English, particularly in terms of script, vocabulary, and grammatical structures (Schepens et al., 2012; Gnanadesikan, 2017). For languages that are structurally closer to English, models tend to perform at levels more comparable to their English proficiency (Conneau, 2019; Pires, 2019), which makes it more challenging to uncover their cross-lingual weaknesses. A more detailed analysis

Table 1: Comparison of conversion rates across different languages. **NP** (No Perturbation) refers to direct translation without perturbations, while **DP** (Direct Perturbation) applies perturbations without search.

Language	NP	DP	Ours
Chinese	0.000	0.036	0.431
Japanese	0.000	0.071	0.594
French	0.000	0.018	0.132
German	0.000	0.027	0.323

of how linguistic relationships affect cross-lingual model performance is provided in subsection 3.3.

**Our search framework significantly outperforms baseline approaches.** We compared our beam search method with two baseline approaches: No Perturbation (NP) and Direct Perturbation (DP). In NP, we directly translate the English questions to target languages without any perturbation, while in DP, we apply perturbations through prompts following the template in Appendix E directly, without search. Using models in  $\mathcal{M}$  for simulation, we identify questions where models perform well in English but fail in target languages. As shown in Table 1, our framework consistently achieves substantially higher conversion rates across all evaluated languages compared to both baselines.

### 3.3 Linguistic Factors in Cross-lingual Weaknesses

To answer RQ2, we first sampled 100 seed bilingual pairs (i.e., English-target language pairs) for each of 16 languages from those generated in subsection 3.2. For each sampled pair, the target-language portion was translated into the other 15 languages, resulting in a total of 25,600 expanded bilingual pairs. These expanded pairs were then evaluated across six different models, with detailed experimental settings outlined in Appendix B.

The identified cross-lingual weaknesses are not restricted to specific languages and depend on the linguistic relationships. The evaluation results of GPT-4o-mini on expanded pairs from the Asian language family (Chinese, Japanese, and Korean) and the European language family (French, German, and Spanish) are presented in Figure 6. As observed, the model exhibits a consistent decline in accuracy across these pairs.

A clear pattern emerges when analyzing the expanded pairs. Within the Asian language family, weakness pairs expanded from Chinese, Japanese, or Korean into other Asian languages exhibit substantial and relatively consistent accuracy declines. In contrast, when these Asian seed pairs are expanded into European languages, the accuracy drops are considerably smaller and more variable. A similar trend is observed within the European language family: pairs expanded from French, German, or Spanish into other European languages experience significant and consistent accuracy declines, whereas expansion into Asian languages results in smaller and more varied reductions in accuracy. We hypothesize that these patterns are driven by underlying linguistic relationships. Specifically, the Asian language family exhibits shared cross-linguistic challenges, while the European family follows similar patterns. Consequently, expanded pairs from Chinese seed pairs tend to maintain more weakness in Japanese and Korean, whereas those from French seed pairs lead to increased weakness in German and Italian.

**Languages with stronger linguistic affinity tend to exhibit cross-lingual weaknesses in common.** We define the **Relative Affinity Score (RAS)**  $D_{x,y}$ , which measures the linguistic relationships between language  $x$  and language  $y$ . The score is computed as:

$$D_{x,y} = \left( \frac{A_{x,y} - \bar{A}_x}{\bar{A}_x} \right) \cdot \exp(c \cdot |\bar{A}_y - \bar{A}_x|)$$

Here,  $D_{x,y}$  quantifies the linguistic proximity between language  $x$  and  $y$ , with lower values indicating a stronger affinity. The term  $A_{x,y}$  represents the model’s accuracy on language  $x$  when using expanded pairs originating from seed language  $y$ . The average accuracy on language  $x$  across all seed languages for expanded pairs is denoted by  $\bar{A}_x$ . The factor  $\exp(c \cdot |\bar{A}_y - \bar{A}_x|)$  scales the score based on the accuracy difference between languages  $x$  and  $y$ , where constant  $c$ , a negative value, controls the inverse sensitivity of this adjustment.

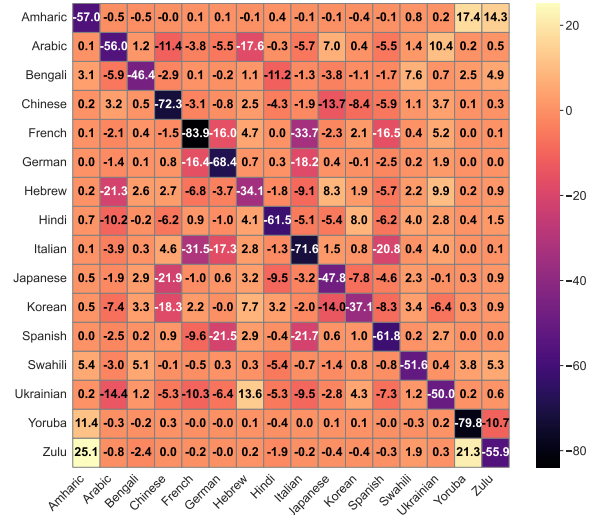


Figure 7: Visualization of RAS  $D_{x,y}$  across 16 languages, highlighting linguistic and cultural proximities. The vertical axis denotes the source language and the horizontal axis denotes the target language. Darker shades of a block indicate stronger retention of shared cross-lingual weaknesses when translating from language  $y$  to language  $x$ , signifying a closer linguistic relationship between the two languages.

As shown in Figure 7, it reveals a clear pattern: lower RAS values  $D_{x,y}$  are predominantly observed for language pairs  $(x, y)$  with linguistic and cultural proximities. This observation strongly supports our hypothesis that languages with closer linguistic ties tend to share cross-lingual weaknesses.

We further investigated the linguistic basis of these cross-lingual weaknesses by analyzing the embedding space of seed bilingual pairs. Specifically, we embedded cross-lingual weaknesses identified in subsection 3.2 for the Asian and European language families. As shown in Figure 9, visualizing these embeddings via t-SNE revealed a clustering effect: weaknesses from the same family clustered together. This observation was corroborated by the cosine distance matrix, as presented in Figure 8, which showed significantly smaller embedding distances within the Asian and European families compared to distances between families.

**Cross-lingual weaknesses correlate with specific subject domains.** To further investigate cross-lingual weaknesses, we categorized the identified bilingual pairs into six subject domains: Science & Technology, History & World Events, Society & Culture, Arts & Literature, Geography & Environment, and General Knowledge. The distribution of these weaknesses across categories, broken down by language, is detailed in Table 4. Notably, lower-resource languages such as Amharic, Arabic, and

Table 2: Performance comparison of Phi-3.5-Mini and Llama-3.1-8B after SFT and DPO on French and Chinese datasets. The table shows evaluation results on various evaluation languages (EL), with the Asian language group highlighted in blue. Performance differences (Diff.) are shown compared to the original model (Orig.). "S Enh." represents the model enhanced by SFT, and "D Enh." represents the model enhanced by simulated DPO. Due to space limitations, performance of Gemma-2-9B and Qwen2.5-7B are presented in Table 5.

Model	EL	Orig.	French Fine-Tuning				Chinese Fine-Tuning			
			S Enh.	D Enh.	S Diff.	D Diff.	S Enh.	D Enh.	S Diff.	D Diff.
Phi-3.5-Mini	French	0.196	–	–	–	–	0.397	0.425	0.202	0.229
	German	0.208	0.466	0.488	0.258	0.258	0.441	0.465	0.233	0.257
	Spanish	0.248	0.467	0.495	0.219	0.247	0.521	0.550	0.273	0.302
	Chinese	0.199	0.345	0.330	0.146	0.131	–	–	–	–
	Japanese	0.127	0.306	0.325	0.178	0.198	0.478	0.510	0.350	0.383
	Korean	0.086	0.246	0.265	0.160	0.179	0.458	0.492	0.373	0.406
Llama-3.1-8B	French	–	–	–	–	–	0.506	0.518	0.343	0.355
	German	0.189	0.388	0.395	0.199	0.206	0.494	0.502	0.304	0.312
	Spanish	0.145	0.409	0.418	0.264	0.273	0.438	0.445	0.293	0.300
	Chinese	0.289	0.453	0.445	0.164	0.156	–	–	–	–
	Japanese	0.124	0.357	0.368	0.232	0.243	0.561	0.570	0.436	0.445
	Korean	0.143	0.320	0.328	0.178	0.185	0.504	0.515	0.362	0.373

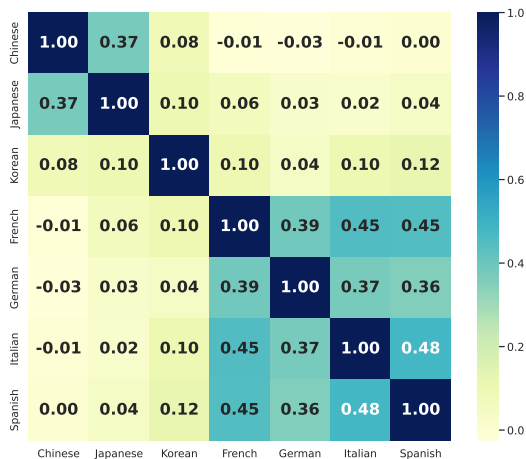


Figure 8: Heatmap of the pairwise cosine distances between the normalized embeddings generated by Llama-3.1-8B for seven English–target language question pairs.

Yoruba exhibited significantly more errors in the Science & Technology domain compared to most languages. Conversely, higher-resource languages like Chinese, Spanish, and German demonstrated stronger performance in this area. Interestingly, Chinese showed a distinct weakness in Society & Culture, while Korean displayed comparatively weaker performance in Geography & Environment.

### 3.4 Cross-lingual Fine-tuning

To answer RQ3, we designed a fine-tuning experiment to explore whether language-specific fine-tuning preferentially enhances performance on linguistically similar languages. We focused on two language families identified as linguistically proximate in our earlier analysis: the Asian language family (Chinese, Japanese, and Korean) and the

European language family (French, German, and Spanish). Using the Chinese and French seed pairs identified in subsection 3.2, we performed both supervised fine-tuning (SFT) and Direct Preference Optimization (DPO) (Rafailov et al., 2023) on several LLMs: Phi-3.5-Mini, Gemma-2-9B, Llama-3.1-8B, Qwen2.5-7B. Separate fine-tuning runs were conducted using both Chinese and French portion in seed pairs. Subsequently, we evaluated the performance of these fine-tuned models across the other languages in two language families. This experiment aimed to investigate if fine-tuning on a specific language leads to greater performance gains in linguistically related languages.

As shown in Table 2 and Table 5, the evaluation results reveal a consistent trend: fine-tuning on Chinese significantly improves performance in Japanese and Korean, while its impact on European languages is comparatively smaller. Similarly, fine-tuning on French enhances performance in related European languages like German and Spanish but has a weaker effect on Asian languages. This pattern holds across both SFT and DPO fine-tuning, indicating that linguistic proximity, rather than the fine-tuning method, primarily drives cross-lingual knowledge transfer. These findings suggest that current LLMs inherently capture linguistic relationships, facilitating more effective transfer between closely related languages.

## 4 Discussion

Our investigation into cross-lingual weaknesses underscores several critical aspects for both under-



standing current LLM limitations and paving the way for future improvements.

First, the integrity of our findings hinges on the **quality of translation** in bilingual question pairs. If semantic equivalence between the English source and target language question is not rigorously maintained, observed performance drops could be mistakenly attributed to the model’s cross-lingual deficiencies rather than translation artifacts. To mitigate this, we employed LLMs for both initial translation and semantic verification, a widely adopted practice in multilingual research (Lin et al., 2024; Ye et al., 2024). The efficacy of this approach was further corroborated through human evaluation, whose methodology and results are presented in Appendix C. The evaluation confirmed that most generated pairs exhibit high translational fidelity. As multilingual capabilities of LLMs continue to advance, developing more sophisticated and reliable translation and semantic checking components will be instrumental in refining the precision with which cross-lingual weaknesses are identified and analyzed.

Second, to provide a richer, more nuanced understanding beyond aggregate statistics, we have compiled an extensive set of **case studies**. These qualitative examples, detailed in Appendix D, illustrate the diverse nature of cross-lingual pitfalls encountered by various models across different languages. They showcase specific failure modes, such as misinterpretation of nuanced phrasing, incorrect entity mapping, or breakdowns in reasoning when faced with linguistic structures that differ significantly from English. These case studies offer valuable material for researchers seeking to conduct in-depth analyses of specific cross-lingual phenomena or to understand the particular challenges faced by individual models or language families.

Finally, the identification of these cross-lingual weaknesses is not merely an academic exercise but offers substantial **potential for enhancing the multilingual capabilities of LLMs**. Our methodology serves as a diagnostic tool, pinpointing specific areas where LLMs falter, thereby guiding targeted interventions. For instance, the weaknesses uncovered can inform more focused **fine-tuning strategies**, concentrating efforts on language pairs or specific linguistic constructions where models demonstrate pronounced deficiencies, potentially leveraging the subject domain categorizations (as shown in Table 4) to further refine this targeting. Furthermore, the challenging cross-lingual exam-

ples generated by our method can be invaluable for **augmenting pre-training and instruction-tuning datasets** (Huang et al., 2024b). By enriching training corpora with instances that expose known weaknesses, we can proactively address data imbalances or representational gaps that contribute to these performance discrepancies. Lastly, these targeted examples are well-suited for **continual learning or adaptive training paradigms**, enabling models to iteratively strengthen their cross-lingual understanding and reasoning in precisely the areas where they have been shown to be vulnerable. In essence, a systematic approach to uncovering weaknesses, such as the one proposed, is a crucial first step towards building more robust multilingual LLMs.

## 5 Conclusion

In this study, we proposed an efficient beam search with LLM-based simulation to identify cross-lingual weaknesses in LLMs, generating a 16-language dataset that exposed performance gaps even in state-of-the-art models. Our findings highlight linguistic relationships as key to shared vulnerabilities and fine-tuning benefits, emphasizing the need to consider linguistic nuances in developing truly multilingual LLMs.

## Limitations

While our methodology demonstrates effectiveness in identifying cross-lingual weaknesses, several avenues for future refinement exist. First, the current study’s scope, while covering a diverse set of languages, is not fully comprehensive. A more complete picture of cross-lingual consistency in LLMs would require extending our analysis to a broader range of languages, particularly those with limited resources or significantly different structural characteristics. Relatedly, although we employ LLM-based semantic checks to ensure the semantic equivalence of our bilingual question pairs, subtle nuances arising from cultural context or idiomatic expressions might still introduce minor biases. Finally, our core approach of iteratively adding perturbations is effective at revealing weaknesses related to complexity. However, this strategy may be less sensitive to identifying those vulnerabilities that manifest in very short, concise prompts. Consequently, investigating complementary techniques specifically designed for such cases would enhance the overall robustness of our framework.

## Ethics Statement

This research adheres to ethical standards in AI research and development. Our methodology is designed to identify and understand cross-lingual weaknesses in LLMs to improve their multilingual capabilities. We recognize the potential for bias within LLMs, particularly across different languages and cultural contexts. Our language selection was carefully considered to ensure diversity, encompassing both high-resource and lower-resource languages. All generated content and model outputs were scrutinized for potential biases. No personally identifiable information was collected or used. This work is intended to promote inclusivity and fairness in the development of multilingual LLMs. The findings are shared with the research community to foster further investigation and the mitigation of cross-lingual weaknesses in LLMs.

## References

- Marah Abidin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Anthropic. 2024. Claude 3.5 sonnet. <https://www.anthropic.com/news/claude-3-5-sonnet>.
- Prajwal Bhargava, Aleksandr Drozd, and Anna Rogers. 2021. Generalization in nli: Ways (not) to go beyond simple heuristics. *arXiv preprint arXiv:2110.01518*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Xiuying Chen, Tairan Wang, Taicheng Guo, Kehan Guo, Juexiao Zhou, Haoyang Li, Zirui Song, Xin Gao, and Xiangliang Zhang. 2025. Unveiling the power of language models in chemical research question answering. *Communications Chemistry*, 8(1):4.
- Gao Chujie, Siyuan Wu, Yue Huang, Dongping Chen, Qihui Zhang, Zhengyan Fu, Yao Wan, Lichao Sun, and Xiangliang Zhang. 2024. Honestllm: Toward an honest and helpful large language model. *Advances in Neural Information Processing Systems*, 37:7213–7255.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.
- A Conneau. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.
- Sumanth Doddapaneni, Mohammed Safi Ur Rahman Khan, Dilip Venkatesh, Raj Dabre, Anoop Kunchukuttan, and Mitesh M Khapra. 2024. Cross-lingual auto evaluation for assessing multilingual llms. *arXiv preprint arXiv:2410.13394*.
- Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Hiroki Iida, Masanari Ohi, Kakeru Hattori, Hirai Shota, Sakae Mizuki, Rio Yokota, and Naoaki Okazaki. 2024. Continual pre-training for cross-lingual llm adaptation: Enhancing japanese language capabilities. *arXiv preprint arXiv:2404.17790*.
- Lang Gao, Xiangliang Zhang, Preslav Nakov, and Xiuying Chen. 2024. [Shaping the safety boundaries: Understanding and defending against jailbreaks in large language models](#). *Preprint*, arXiv:2412.17034.
- Amalia E Gnanadesikan. 2017. Towards a typology of phonemic scripts. *Writing Systems Research*, 9(1):14–35.
- Google. [Google translate api](#). Accessed: 2025-01-11.
- Taicheng Guo, Bozhao Nan, Zhenwen Liang, Zhichun Guo, Nitesh Chawla, Olaf Wiest, Xiangliang Zhang, et al. 2023. What can large language models do in chemistry? a comprehensive benchmark on eight tasks. *Advances in Neural Information Processing Systems*, 36:59662–59688.
- Katharina Hämmerl, Jindřich Libovický, and Alexander Fraser. 2024. Understanding cross-lingual alignment—a survey. *arXiv preprint arXiv:2404.06228*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2021a. Aligning ai with shared human values. *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021b. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021c. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Yue Huang, Chenrui Fan, Yuan Li, Siyuan Wu, Tianyi Zhou, Xiangliang Zhang, and Lichao Sun. 2024a. 1+ 1 > 2: Can large language models serve as cross-lingual knowledge aggregators? *arXiv preprint arXiv:2406.14721*.

- Yue Huang, Chujie Gao, Siyuan Wu, Haoran Wang, Xiangqi Wang, Yujun Zhou, Yanbo Wang, Jiayi Ye, Jiawen Shi, Qihui Zhang, et al. 2025a. On the trustworthiness of generative foundation models: Guideline, assessment, and perspective. *arXiv preprint arXiv:2502.14296*.
- Yue Huang, Yanbo Wang, Zixiang Xu, Chujie Gao, Siyuan Wu, Jiayi Ye, Xiuying Chen, Pin-Yu Chen, and Xiangliang Zhang. 2025b. Breaking focus: Contextual distraction curse in large language models. *arXiv preprint arXiv:2502.01609*.
- Yue Huang, Siyuan Wu, Chujie Gao, Dongping Chen, Qihui Zhang, Yao Wan, Tianyi Zhou, Chaowei Xiao, Jianfeng Gao, Lichao Sun, et al. 2024b. Datagen: Unified synthetic dataset generation via large language models. In *The Thirteenth International Conference on Learning Representations*.
- Yue Huang, Zhengqing Yuan, Yujun Zhou, Kehan Guo, Xiangqi Wang, Haomin Zhuang, Weixiang Sun, Lichao Sun, Jindong Wang, Yanfang Ye, et al. 2024c. Social science meets llms: How reliable are large language models in social simulations? *arXiv preprint arXiv:2410.23426*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. 2024. Openai o1 system card. *arXiv preprint arXiv:2412.16720*.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*.
- Yuan Li, Yue Huang, Hongyi Wang, Xiangliang Zhang, James Zou, and Lichao Sun. 2024a. Quantifying ai psychology: A psychometrics benchmark for large language models. *arXiv preprint arXiv:2406.17675*.
- Zihao Li, Yucheng Shi, Zirui Liu, Fan Yang, Ali Payani, Ninghao Liu, and Mengnan Du. 2024b. **Language ranker: A metric for quantifying llm performance across high and low-resource languages**. *Preprint*, arXiv:2404.11553.
- Luyang Lin, Lingzhi Wang, Jinsong Guo, and Kam-Fai Wong. 2024. Investigating bias in llm-based bias detection: Disparities between llms and human perception. *arXiv preprint arXiv:2403.14896*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*.
- Hongwei Liu, Zilong Zheng, Yuxuan Qiao, Haodong Duan, Zhiwei Fei, Fengzhe Zhou, Wenwei Zhang, Songyang Zhang, Dahua Lin, and Kai Chen. 2024a. Mathbench: Evaluating the theory and application proficiency of llms with a hierarchical mathematics benchmark. *arXiv preprint arXiv:2405.12209*.
- Wei Liu, Zhongyu Niu, Lang Gao, Zhiying Deng, Jun Wang, Haozhao Wang, and Ruixuan Li. 2025. **Adversarial cooperative rationalization: The risk of spurious correlations in even clean datasets**. *Preprint*, arXiv:2505.02118.
- Wei Liu, Chenxi Wang, YiFei Wang, Zihao Xie, Rennai Qiu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, and Chen Qian. 2024b. **Autonomous agents for collaborative task under information asymmetry**. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Zhuoer Feng, Bosi Wen, Jiale Cheng, Pei Ke, Yifan Xu, Weng Lam Tam, et al. 2023a. Alignbench: Benchmarking chinese alignment of large language models. *arXiv preprint arXiv:2311.18743*.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. 2023b. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*.
- Yu Liu, Lang Gao, Mingxin Yang, Yu Xie, Ping Chen, Xiaojin Zhang, and Wei Chen. 2024c. **Vuldetbench: Evaluating the deep capability of vulnerability detection with large language models**. *Preprint*, arXiv:2406.07595.
- Haoran Luo, E Haihong, Yuhao Yang, Gengxian Zhou, Yikai Guo, Tianyu Yao, Zichen Tang, Xueyuan Lin, and Kaiyang Wan. 2023. Nqe: N-ary query embedding for complex query answering over hyper-relational knowledge graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 4543–4551.
- Meta. 2024a. Llama 3.1-70b. <https://huggingface.co/meta-llama/Llama-3.1-70B>.
- Meta. 2024b. Llama 3.1-8b. <https://huggingface.co/meta-llama/Llama-3.1-8B>.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.
- Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. 2023. A comprehensive overview of large language models. *arXiv preprint arXiv:2307.06435*.

- OpenAI. 2024. Gpt-4o mini: Advancing cost-efficient intelligence. <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>.
- T Pires. 2019. How multilingual is multilingual bert. *arXiv preprint arXiv:1906.01502*.
- Libo Qin, Qiguang Chen, Fuxuan Wei, Shijue Huang, and Wanxiang Che. 2023. Cross-lingual prompting: Improving zero-shot chain-of-thought reasoning across languages. *arXiv preprint arXiv:2310.14799*.
- 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.
- Job Schepens, Ton Dijkstra, and Franc Grootjen. 2012. Distributions of cognates in europe as based on levenshtein distance. *Bilingualism: Language and Cognition*, 15(1):157–166.
- Zirui Song, Bin Yan, Yuhan Liu, Miao Fang, Mingzhe Li, Rui Yan, and Xiuying Chen. 2025a. Injecting domain-specific knowledge into large language models: a comprehensive survey. *arXiv preprint arXiv:2502.10708*.
- Zirui Song, Jingpu Yang, Yuan Huang, Jonathan Tonglet, Zeyu Zhang, Tao Cheng, Meng Fang, Iryna Gurevych, and Xiuying Chen. 2025b. Geolocation with real human gameplay data: A large-scale dataset and human-like reasoning framework. *arXiv preprint arXiv:2502.13759*.
- Joe Stacey, Pasquale Minervini, Haim Dubossarsky, Sebastian Riedel, and Tim Rocktäschel. 2020. Avoiding the hypothesis-only bias in natural language inference via ensemble adversarial training. *arXiv preprint arXiv:2004.07790*.
- Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. 2024. Trustllm: Trustworthiness in large language models. *arXiv preprint arXiv:2401.05561*, 3.
- Gemma Team. 2024a. *Gemma*.
- Qwen Team. 2024b. *Qwen2.5: A party of foundation models*.
- Alan Wake, Albert Wang, Bei Chen, CX Lv, Chao Li, Chengen Huang, Chenglin Cai, Chujie Zheng, Daniel Cooper, Ethan Dai, et al. 2024. Yi-lightning technical report. *arXiv e-prints*, pages arXiv–2412.
- Kaiyang Wan, Honglin Mu, Rui Hao, Haoran Luo, Tianle Gu, and Xiuying Chen. 2025. A cognitive writing perspective for constrained long-form text generation. *arXiv preprint arXiv:2502.12568*.
- Chenxi Wang, Tianle Gu, Zhongyu Wei, Lang Gao, Zirui Song, and Xiuying Chen. 2025a. Word form matters: Llms’ semantic reconstruction under typoglycemia. *Preprint*, arXiv:2503.01714.
- Chenxi Wang, Zongfang Liu, Dequan Yang, and Xiuying Chen. 2025b. Decoding echo chambers: LLM-powered simulations revealing polarization in social networks. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 3913–3923, Abu Dhabi, UAE. Association for Computational Linguistics.
- Hetong Wang, Pasquale Minervini, and Edoardo M Ponti. 2024. Probing the emergence of cross-lingual alignment during llm training. *arXiv preprint arXiv:2406.13229*.
- Jiaan Wang, Yunlong Liang, Zengkui Sun, Yuxuan Cao, Jiarong Xu, and Fandong Meng. 2023. Cross-lingual knowledge editing in large language models. *arXiv preprint arXiv:2309.08952*.
- Yanbo Wang, Jiayi Ye, Siyuan Wu, Chujie Gao, Yue Huang, Xiuying Chen, Yue Zhao, and Xiangliang Zhang. 2025c. Trusteval: A dynamic evaluation toolkit on trustworthiness of generative foundation models. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (System Demonstrations)*, pages 70–84.
- Yumeng Wang, Zhiyuan Fan, Qingyun Wang, May Fung, and Heng Ji. 2025d. Calm: Unleashing the cross-lingual self-aligning ability of language model question answering. *arXiv preprint arXiv:2501.18457*.
- Johannes Welbl, Nelson F Liu, and Matt Gardner. 2017. Crowdsourcing multiple choice science questions. *arXiv preprint arXiv:1707.06209*.
- Yunfei Xie, Ce Zhou, Lang Gao, Juncheng Wu, Xianhang Li, Hong-Yu Zhou, Sheng Liu, Lei Xing, James Zou, Cihang Xie, and Yuyin Zhou. 2025. Medtrinity-25m: A large-scale multimodal dataset with multigranular annotations for medicine. *Preprint*, arXiv:2408.02900.
- Rongwu Xu, Zehan Qi, Zhijiang Guo, Cunxiang Wang, Hongru Wang, Yue Zhang, and Wei Xu. 2024. Knowledge conflicts for llms: A survey. *arXiv preprint arXiv:2403.08319*.
- Zixiang Xu, Yanbo Wang, Yue Huang, Jiayi Ye, Haomin Zhuang, Zirui Song, Lang Gao, Chenxi Wang, Zhaorun Chen, Yujun Zhou, Sixian Li, Wang Pan, Yue Zhao, Jieyu Zhao, Xiangliang Zhang, and Xiuying Chen. 2025a. Socialmaze: A benchmark for evaluating social reasoning in large language models.
- Zixiang Xu, Yanbo Wang, Chenxi Wang, Lang Gao, Zirui Song, Yue Huang, Zhaorun Chen, Xiangliang Zhang, and Xiuying Chen. 2025b. Gta: Graph theory agent and benchmark for algorithmic graph reasoning with llms.

- Ikuya Yamada and Ryokan Ri. 2024. Leia: Facilitating cross-lingual knowledge transfer in language models with entity-based data augmentation. *arXiv preprint arXiv:2402.11485*.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, and Others. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
- Binwei Yao, Ming Jiang, Diyi Yang, and Junjie Hu. 2023. Benchmarking llm-based machine translation on cultural awareness. *arXiv preprint arXiv:2305.14328*.
- Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, et al. 2024. Justice or prejudice? quantifying biases in llm-as-a-judge. *arXiv preprint arXiv:2410.02736*.
- Biao Zhang, Barry Haddow, and Alexandra Birch. 2023a. Prompting large language model for machine translation: A case study. In *International Conference on Machine Learning*, pages 41092–41110. PMLR.
- Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. 2023b. Sentiment analysis in the era of large language models: A reality check. *arXiv preprint arXiv:2305.15005*.
- Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji Kawaguchi, and Lidong Bing. 2024a. How do large language models handle multilingualism? *arXiv preprint arXiv:2402.18815*.
- Yiran Zhao, Wenxuan Zhang, Huiming Wang, Kenji Kawaguchi, and Lidong Bing. 2024b. Adamergex: Cross-lingual transfer with large language models via adaptive adapter merging. *arXiv preprint arXiv:2402.18913*.
- Juexiao Zhou, Xiuying Chen, and Xin Gao. 2023. Path to medical agi: Unify domain-specific medical llms with the lowest cost. *arXiv preprint arXiv:2306.10765*.

## A Related Work

### A.1 LLM Evaluation

Significant efforts have been devoted to evaluating the capabilities of LLMs across a wide range of domains. These evaluations include traditional NLP tasks such as sentiment analysis (Zhang et al., 2023b; Wan et al., 2025) and translation (Yao et al., 2023; Zhang et al., 2023a), as well as mathematical reasoning (Hendrycks et al., 2021c; Liu et al., 2024a), scientific and domain-specific question answering (Xu et al., 2025b; Luo et al., 2023; Zhou et al., 2023; Chen et al., 2025; Song et al., 2025a), and coding skills (Chen et al., 2021; Jain et al., 2024). Evaluations have also extended into specialized domains such as chemistry (Chen et al., 2025), medicine (Xie et al., 2025), and geolocation reasoning (Song et al., 2025b). In the area of cybersecurity, efforts have been made to assess LLMs’ ability to detect software vulnerabilities (Liu et al., 2024c). Beyond task performance, growing attention has been paid to trustworthiness (Sun et al., 2024; Chujie et al., 2024), including robustness to spurious correlations (Liu et al., 2025), resilience to textual perturbations (Wang et al., 2025a), and defense against jailbreak attacks (Gao et al., 2024). Comprehensive benchmarks and investigations have been proposed to systematically assess these aspects (Huang et al., 2025a; Wang et al., 2025c). General-purpose benchmarks like MMLU (Hendrycks et al., 2021b,a) continue to serve as a foundation for evaluating broad LLM capabilities.

In this study, we select a subset of English questions from five widely used question-answering datasets: CommonsenseQA (Naveed et al., 2023), ARC (Clark et al., 2018), MMLU (Hendrycks et al., 2021b,a), SciQ (Welbl et al., 2017), and TruthfulQA (Lin et al., 2021). These datasets evaluate models on common sense reasoning, mathematical problem-solving, scientific knowledge, and various other skills. We use these questions as the foundation for generating our own dataset.

### A.2 Cross-lingual Capability of LLMs.

The cross-lingual capabilities of LLMs have become a central focus in NLP research. Multi-task finetuning (MTF) has proven effective for enhancing cross-lingual generalization, as shown by Muennighoff et al. (2022), where finetuning multilingual models like BLOOM and mT5 on English tasks enabled zero-shot task transfer to other

languages. Beyond MTF, cross-lingual prompting techniques such as chain-of-thought (CoT) prompting (Qin et al., 2023) improve reasoning accuracy by aligning representations and employing task-specific solvers. Other approaches, including cross-lingual knowledge editing (Wang et al., 2023), entity-based data augmentation (Yamada and Ri, 2024) and cross-lingual knowledge aggregator (Huang et al., 2024a), have been proposed to enhance adaptation and infuse models with cross-lingual knowledge.

Evaluation has also gained attention, with frameworks like the Cross Lingual Auto Evaluation (CIA) Suite (Doddapaneni et al., 2024) addressing challenges in assessing multilingual model outputs. However, many MTF studies remain English-centric (Muennighoff et al., 2022), and prompting techniques (Qin et al., 2023) may struggle with diverse linguistic structures. While methods like adapter merging (Zhao et al., 2024b) and continual pre-training (Fujii et al., 2024) aim to enhance language transfer, systematic investigation into multilingual LLM weaknesses across diverse languages remains limited. Additionally, while studies probe cross-lingual alignment during pre-training (Wang et al., 2024; Liu et al., 2023a) and its importance (Hämmerl et al., 2024), a quantifiable measure of linguistic relationships affecting cross-lingual transfer is absent.

Our work builds on these foundations by systematically identifying and analyzing cross-lingual weaknesses in LLMs across 16 diverse languages. By introducing a novel metric to quantify linguistic relationships based on observed performance, we offer deeper insights into how linguistic relation impacts model behavior.

## B Experiment Details

### B.1 Experiment Settings

**Source dataset.** To create bilingual pairs, we randomly sampled English questions from five commonly used datasets that cover a wide range of model capabilities: ARC, MMLU, CommonsenseQA, TruthfulQA, and SciQ. The sampling was performed equally across all five datasets.

**Models.** As detailed in Table 3, we utilized five proprietary models: GPT-4o (Hurst et al., 2024), GPT-4o-mini (OpenAI, 2024), Yi-Lightning (Wake et al., 2024), Claude-3.5-Sonnet (Anthropic, 2024), and o1-mini (Jaech et al., 2024). In addition, we included seven open-weight models: Gemma-

Table 3: Models used in our experiments along with their versions, organizations, licenses, and purposes. *Eval*: Model used for evaluation; *FT*: Model used for fine-tuning.

Model	Version	Organization	License	Eval	FT
GPT-4o-mini	gpt-4o-mini-2024-07-18	OpenAI	Proprietary	✓	
GPT-4o	gpt-4o-2024-08-06	OpenAI	Proprietary	✓	
Gemma-2-9B	Gemma-2-9B-it	Google	Gemma License	✓	✓
Gemma-2-27B	Gemma-2-27B-it	Google	Gemma License	✓	
Llama-3.1-8B	Meta-Llama-3.1-8B-Instruct	Meta	Llama 3.1 Community	✓	✓
Llama-3.1-70B	Meta-Llama-3.1-70B-Instruct	Meta	Llama 3.1 Community	✓	
Yi-Lightning	Yi-Lightning	01 AI	Proprietary	✓	
Qwen2.5-7B	Qwen2.5-7B-Instruct	Alibaba	Qwen License	✓	✓
Qwen2.5-72B	Qwen2.5-72B-Instruct	Alibaba	Qwen License	✓	
o1-mini	o1-mini-2024-09-12	OpenAI	Proprietary	✓	
Phi-3.5-mini	Phi-3.5-mini-instruct	Microsoft	MIT		✓
Claude-3.5-Sonnet	claude-3-5-sonnet-20241022	Anthropic	Proprietary	✓	

2-9B, Gemma-2-27B (Team, 2024a), Qwen2.5-7B, Qwen2.5-72B (Yang et al., 2024; Team, 2024b), Llama-3.1-8B (Meta, 2024b), Llama-3.1-70B (Meta, 2024a), and Phi-3.5-mini (Abdin et al., 2024).

**Hyperparameter settings.** For perturbation generation, we used a temperature of 0.7 to encourage more diverse and creative responses. In the translation, semantic checking, and simulation tasks, the temperature was reduced to 0.001 to ensure stability in the responses. The maximum output length for these tasks was capped at 1,024 tokens. During beam search, we initialized the process with  $W = 4$  bilingual pairs, and the search width was set to  $w = 12$ . The search depths were configured to  $d_1 = 4$  and  $d_2 = 6$ , respectively. To promote diversity in the generated questions, we set  $r = 3$ . The simulation score parameter,  $\gamma$ , was set to 2. For the Early Stopping Mechanism,  $\theta_{pot}$  was set to 0.6, and for determining inclusion in the candidate list,  $\theta_{inc}$  was set to 0.8. We employed  $K = 5$  LLMs for LLM-based simulation. The constant  $C$  used to calculate the Relative Affinity Score was set to -1.

For the fine-tuning experiments, we trained for 4 epochs with a learning rate of  $2.0e-4$ , employing a cosine learning rate scheduler and a warmup ratio of 0.1. The per-device training batch size was 1, with a gradient accumulation of 8 steps. For evaluation, we used 10% of the training data as a validation set, evaluated every 200 steps, and set the per-device evaluation batch size to 1.

## B.2 Experiment Procedures

**Experiment procedure of cross-Lingual weakness identification.** To generate bilingual question pairs for our cross-lingual weakness identifi-

cation experiments, we employed LLM-based simulation using the following models: Llama-3.1-8B, Gemma-2-9B, Gemma-2-27B, GPT-4o-mini, and Qwen2.5-72B. This process resulted in a total of 6713 bilingual pairs across the following languages: Chinese (342 pairs), Japanese (314 pairs), Korean (456 pairs), French (312 pairs), Spanish (242 pairs), Italian (295 pairs), Ukrainian (323 pairs), German (322 pairs), Bengali (431 pairs), Hindi (327 pairs), Arabic (424 pairs), Hebrew (319 pairs), Amharic (665 pairs), Yoruba (813 pairs), Swahili (417 pairs), and Zulu (711 pairs). Subsequently, we performed zero-shot evaluations on all generated question pairs using the following models: Llama-3.1-8B, Gemma-2-9B, Gemma-2-27B, GPT-4o-mini, Llama-3.1-70B, Qwen2.5-72B, o1-mini, Yi-Lightning, GPT-4o, and Claude-3.5-sonnet. The results of these evaluations are presented in Figure 5, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24 and 25 .

**Experiment procedure of quantifying the linguistic relationships.** To quantify the linguistic relationships between languages, we randomly sampled 100 generated bilingual pairs for each of the following languages: Chinese, Japanese, Korean, French, Spanish, Italian, Ukrainian, German, Bengali, Hindi, Arabic, Hebrew, Amharic, Yoruba, Swahili, and Zulu. We then translated the original question component of these pairs into each of the other fifteen languages using GPT-4o, and the perturbed question component using Google Translate’s API (Google). This process, along with the original language, resulted in a total of 25,600 bilingual pairs (16 languages \* 100 pairs \* 16 translations). We performed zero-shot evaluations on these pairs using six models: Llama-3.1-8B, Gemma-2-27B, GPT-4o-mini, Llama-3.1-70B,

Table 4: Percentage distribution of weaknesses across different categories for each language, compared to overall averages. Percentages exceeding the overall average for each category are highlighted in orange. Column abbreviations are as follows: Sci & Tech (Science & Technology), Gen Knowl. (General Knowledge), Geo & Env. (Geography & Environment), Soc & Cult. (Society & Culture), Arts & Lit. (Arts & Literature), and Hist & Events (History & World Events).

Language	Sci & Tech	Gen Knowl.	Geo & Env.	Soc & Cult.	Arts & Lit.	Hist & Events
Amharic	61.95%	8.12%	5.41%	15.94%	2.26%	6.32%
Arabic	55.42%	15.57%	0.71%	9.20%	9.91%	9.20%
Bengali	46.17%	16.47%	10.21%	17.87%	7.66%	1.62%
Chinese	25.73%	6.43%	7.60%	47.08%	12.28%	0.88%
French	37.50%	13.78%	10.26%	26.28%	6.73%	5.45%
German	42.24%	19.88%	9.32%	22.36%	0.62%	5.59%
Hebrew	43.26%	10.02%	0.58%	22.57%	16.30%	6.27%
Hindi	44.95%	17.74%	3.06%	20.49%	11.31%	2.45%
Italian	49.15%	6.10%	3.39%	26.10%	5.76%	9.49%
Japanese	48.09%	15.29%	12.42%	19.11%	4.46%	0.64%
Korean	27.41%	16.45%	23.25%	25.66%	3.07%	4.17%
Spanish	23.97%	18.18%	4.55%	19.01%	21.90%	12.40%
Swahili	47.96%	8.87%	3.12%	30.94%	4.32%	4.80%
Ukrainian	39.01%	10.53%	4.33%	34.98%	6.81%	4.33%
Yoruba	53.01%	6.52%	7.87%	21.03%	4.67%	6.89%
Zulu	47.40%	13.36%	0.98%	25.60%	8.44%	4.22%
<b>Overall Average</b>	45.36%	12.20%	6.63%	23.40%	7.15%	5.26%

Qwen2.5-72B, and GPT-4o. The Relative Affinity Score was then calculated based on the average accuracy of these models, as shown in Figure 7.

**Experiment procedure of linguistic relationship analysis through fine-tuning.** Leveraging the English-Chinese and English-French question pairs generated in our dataset, we performed SFT and DPO on several Large Language Models: Llama-3.1-8B, Qwen2.5-7B, Gemma-2-9B, and Phi-3.5-Mini. For each model, we conducted separate fine-tuning runs using both the Chinese and French datasets. To ensure consistency across experiments, we trained for 4 epochs with a learning rate of  $2.0e-4$ , employing a cosine learning rate scheduler and a warmup ratio of 0.1. The per-device training batch size was set to 1, with gradient accumulation performed over 8 steps. During training, we used the correct answers from the models’ responses as the target output for each question. For evaluation, we used 10% of the training data as a validation set, evaluated every 200 steps, and set the per-device evaluation batch size to 1.

## C Human Evaluation

To ensure that the target language questions in our generated bilingual pairs maintained semantic equivalence and answer consistency with the original English questions, we conducted a human

evaluation study. We randomly sampled 100 bilingual pairs from the candidate list for each of the following sixteen languages: Chinese, Japanese, Korean, French, Spanish, Italian, Ukrainian, German, Bengali, Hindi, Arabic, Hebrew, Amharic, Yoruba, Swahili, and Zulu. Four undergraduate students majoring in computer science, proficient in English and various translation tools, were divided into two groups to assess: (1) whether the target language question maintained semantic equivalence with the original English question, and (2) whether the answer to the target language question was consistent with the answer to the original English question. The results of this evaluation are summarized in Table 6.

## D Case Study

In Figure 26, 27, 28, 29, 30, 31, 32, and 33, we illustrate case studies of model responses to English-target language (Korean, French, German, Chinese, Italian, Spanish, Japanese, and Ukrainian, respectively) question pairs.



Table 5: Performance comparison of Gemma-2-9B and Qwen2.5-7B after SFT and DPO on French and Chinese datasets. The table shows evaluation results on various evaluation languages (EL), with the Asian language group highlighted in blue. Performance differences (Diff.) are shown compared to the original model (Orig.). "S Enh." represents the model enhanced by SFT, and "D Enh." represents the model enhanced by simulated DPO.

Model	EL	Orig.	French Fine-Tuning				Chinese Fine-Tuning			
			S Enh.	D Enh.	S Diff.	D Diff.	S Enh.	D Enh.	S Diff.	D Diff.
Gemma-2-9B	French	0.222	–	–	–	–	0.510	0.522	0.288	0.300
	German	0.115	0.495	0.505	0.381	0.391	0.463	0.475	0.348	0.360
	Spanish	0.109	0.541	0.555	0.433	0.447	0.463	0.458	0.314	0.309
	Chinese	0.111	0.552	0.560	0.441	0.449	–	–	–	–
	Japanese	0.099	0.527	0.535	0.428	0.436	0.576	0.588	0.478	0.490
	Korean	0.083	0.421	0.430	0.338	0.347	0.537	0.545	0.454	0.462
Qwen2.5-7B	French	0.321	–	–	–	–	0.494	0.503	0.173	0.182
	German	0.233	0.447	0.455	0.214	0.222	0.491	0.485	0.258	0.252
	Spanish	0.145	0.537	0.548	0.393	0.403	0.426	0.432	0.281	0.287
	Chinese	0.281	0.584	0.595	0.304	0.314	–	–	–	–
	Japanese	0.194	0.408	0.415	0.213	0.221	0.592	0.600	0.398	0.406
	Korean	0.140	0.329	0.337	0.189	0.197	0.384	0.393	0.243	0.252

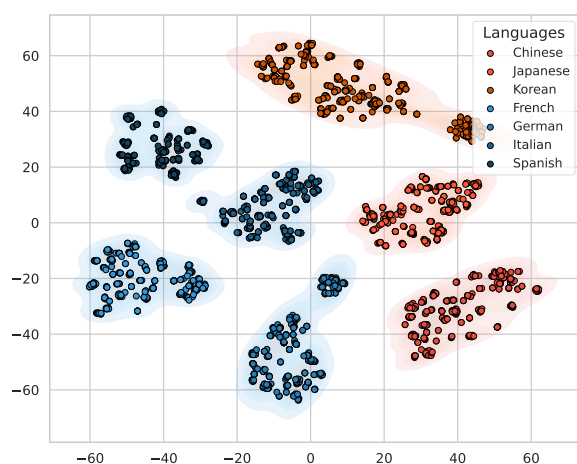


Figure 9: T-SNE visualization of the embeddings generated by LLaMA-3.1-8B for seven English–target language question pairs.

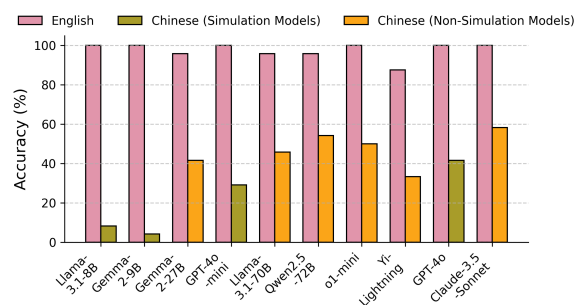


Figure 10: Performance of LLMs on English-Chinese pairs after incorporating GPT-4o into the simulation.

Table 6: Results of human evaluation on semantic equivalence (Semantic Eq.) and answer consistency (Answer Consis.) between original English questions and target language questions in bilingual pairs.

Language	Semantic Eq. (%)	Answer Consis. (%)
<b>Amharic</b>	83.0	88.0
<b>Arabic</b>	90.0	94.0
<b>Bengali</b>	88.0	93.0
<b>Chinese</b>	95.0	98.0
<b>French</b>	97.0	99.0
<b>German</b>	96.0	98.0
<b>Hebrew</b>	93.0	95.0
<b>Hindi</b>	91.0	94.0
<b>Italian</b>	96.0	97.0
<b>Japanese</b>	93.0	95.0
<b>Korean</b>	91.0	93.0
<b>Spanish</b>	98.0	100.0
<b>Swahili</b>	89.0	93.0
<b>Ukrainian</b>	92.0	95.0
<b>Yoruba</b>	84.0	90.0
<b>Zulu</b>	86.0	91.0

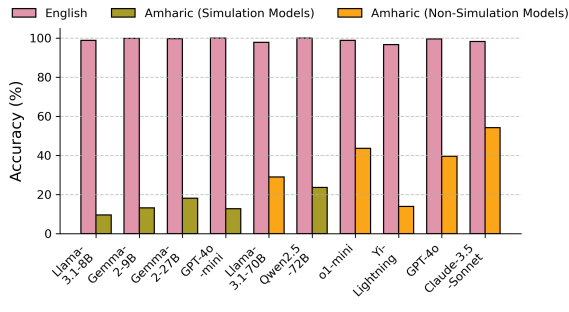


Figure 11: Performance of LLMs on English-Amharic pairs in our candidate list.

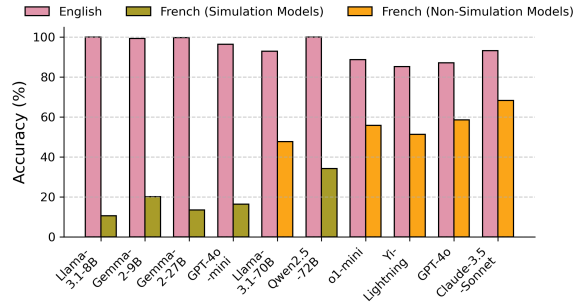


Figure 14: Performance of LLMs on English-French pairs in our candidate list.

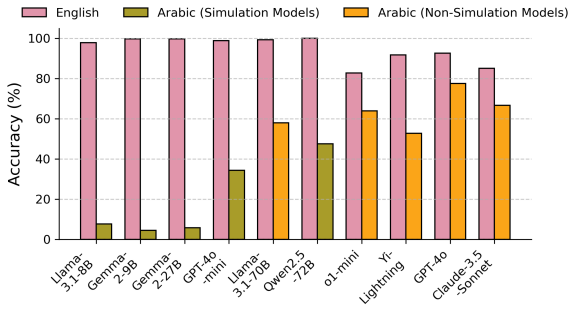


Figure 12: Performance of LLMs on English-Arabic pairs in our candidate list.

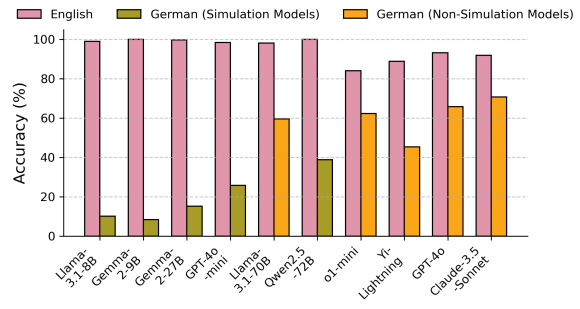


Figure 15: Performance of LLMs on English-German pairs in our candidate list.

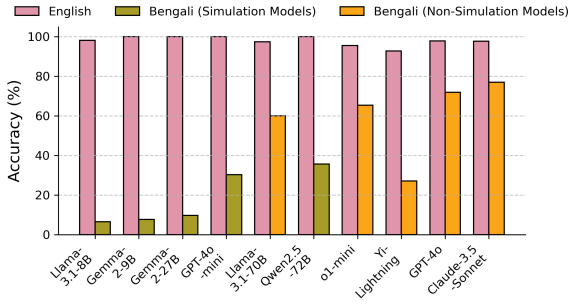


Figure 13: Performance of LLMs on English-Bengali pairs in our candidate list.

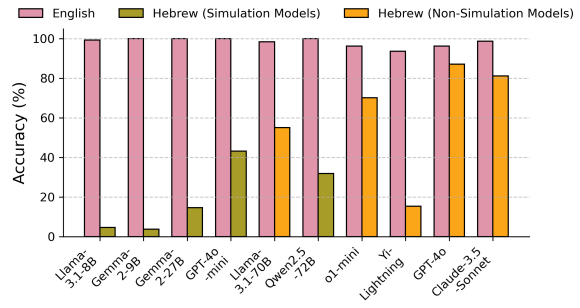


Figure 16: Performance of LLMs on English-Hebrew pairs in our candidate list.

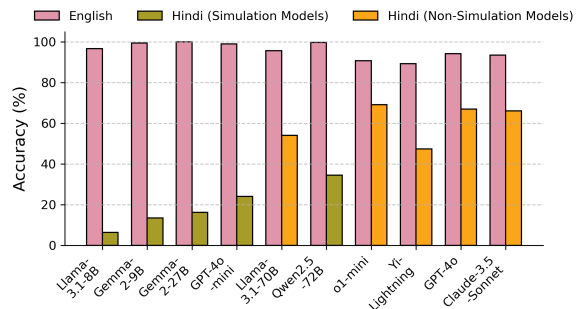


Figure 17: Performance of LLMs on English-Hindi pairs in our candidate list.

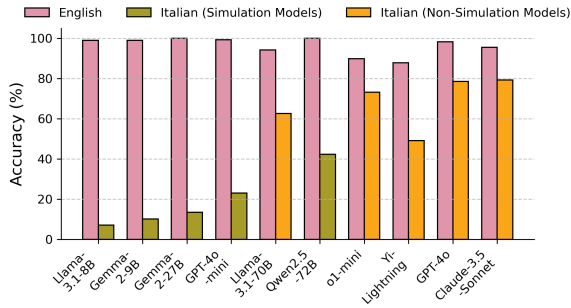


Figure 18: Performance of LLMs on English-Italian pairs in our candidate list.

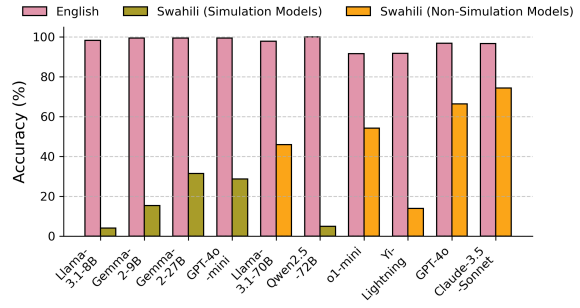


Figure 22: Performance of LLMs on English-Swahili pairs in our candidate list.

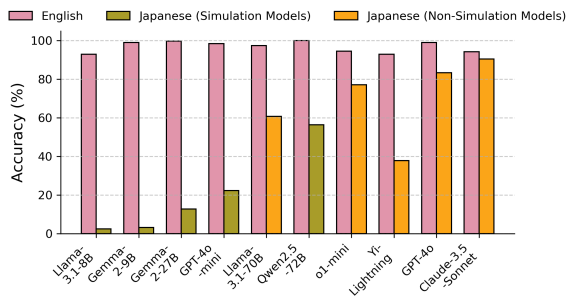


Figure 19: Performance of LLMs on English-Japanese pairs in our candidate list.

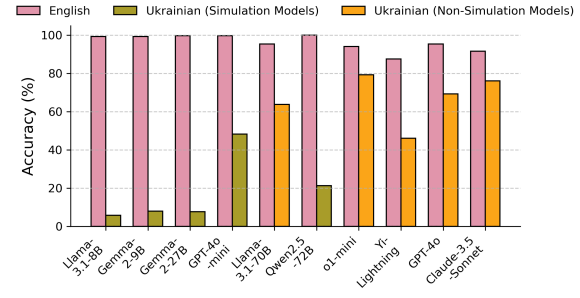


Figure 23: Performance of LLMs on English-Ukrainian pairs in our candidate list.

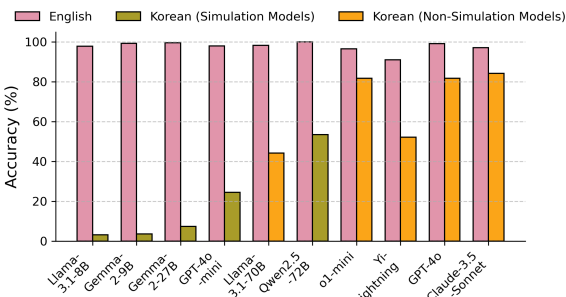


Figure 20: Performance of LLMs on English-Korean pairs in our candidate list.

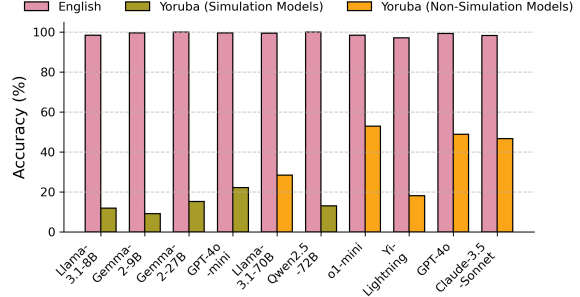


Figure 24: Performance of LLMs on English-Yoruba pairs in our candidate list.

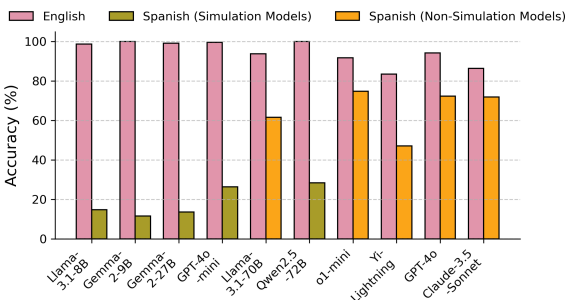


Figure 21: Performance of LLMs on English-Spanish pairs in our candidate list.

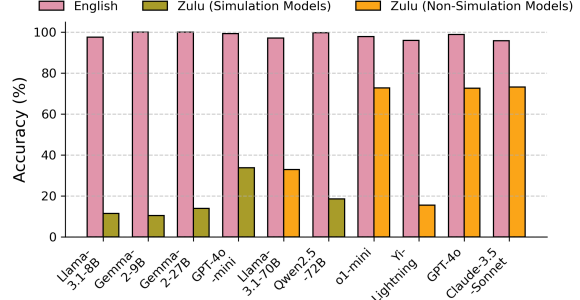


Figure 25: Performance of LLMs on English-Zulu pairs in our candidate list.

**Target Language: Korean**

**English Question:** Many cities around the world, like London, are known for their rich cultural scenes and historic landmarks that attract millions of visitors each year. Kiwis describe it as \"our biggest and most diverse city, it's an economic hub with lots of tourists and very expensive rents\". What is the name of the city? The cultural vibrancy of cities often influences local cuisine, leading to unique dining experiences that reflect their history and diversity.

**English Choices:** [\"New York City\", \"London\", \"Los Angeles\", \"Paris\", \"Toronto\", \"Sydney\", \"Chicago\", \"San Francisco\", \"Wellington\", \"Auckland\"]

**English Ground Truth:** Auckland

---

**Korean Question:** 전 세계의 많은 도시들, 예를 들어 런던은 매년 수백만 명의 방문객을 끌어들이는 풍부한 문화 현장과 역사적인 명소로 유명합니다. 키위들은 이 도시를 '우리의 가장 크고 다양한 도시이며, 경제 중심지로 많은 관광객이 찾고 임대료가 매우 비싸다'라고 설명합니다. 이 도시의 이름은 무엇인가요? 도시의 문화적 활력은 종종 지역 요리에 영향을 미쳐, 그들의 역사와 다양성을 반영하는 독특한 식사 경험을 만들어냅니다.

**Korean Choices:** [\"뉴욕시\", \"런던\", \"로스앤젤레스\", \"파리\", \"토론토\", \"시드니\", \"시카고\", \"샌프란시스코\", \"웰링턴\", \"오클랜드\"]

**Korean Ground Truth:** 오클랜드

**Response by Gemma-2-9B**

<b>English original response:</b> Auckland	<b>Korean original response:</b> 런던
---	--

Figure 26: Case study: Gemma-2-9B’s responses to English-Korean question pairs.

**Target Language: French**

**English Question:** Acute binge drinking is associated with? It's interesting to note how certain lifestyle choices can impact cardiovascular health, often leading to discussions about various heart-related conditions that can be exacerbated by excessive alcohol consumption. This raises an intriguing point about how emotional states can influence physical health, as many studies explore the connection between stress relief and heart function, often highlighting the importance of maintaining a balanced lifestyle for overall well-being.

**English Choices:** ["Happy heart syndrome", "Home heart syndrome", "Beach heart syndrome", "Holiday heart syndrome"]

**English Ground Truth:** Holiday heart syndrome

---

**French Question:** À quoi est associé l'alcoolisation aiguë ? Il est intéressant de noter comment certains choix de mode de vie peuvent impacter la santé cardiovasculaire, menant souvent à des discussions sur diverses affections cardiaques qui peuvent être exacerbées par une consommation excessive d'alcool. Cela soulève un point intrigant sur la façon dont les états émotionnels peuvent influencer la santé physique, car de nombreuses études explorent le lien entre le soulagement du stress et la fonction cardiaque, soulignant souvent l'importance de maintenir un mode de vie équilibré pour le bien-être général.

**French Choices:** ["Syndrome du cœur heureux", "Syndrome du cœur à la maison", "Syndrome du cœur de la plage", "Syndrome du cœur des vacances"]

**French Ground Truth:** Syndrome du cœur des vacances

**Response by GPT-4o-mini**

<p><b>English original response:</b> 'Holiday heart syndrome'</p> <p style="text-align: right;"></p>	<p><b>French original response:</b> 'Syndrome du cœur heureux'</p> <p style="text-align: right;"></p>
---	--

Figure 27: Case study: GPT-4o-mini's responses to English-French question pairs.

**Target Language: German**

**English Question:** Illustrating how form follows function, long, slender protein strands that make up what tissue are essential for contracting and relaxing? The design of various biological structures often reflects their specific roles; for instance, the way some connective tissues are structured can greatly influence their ability to withstand tension and support movement in the body. This principle can also be seen in how tendons connect muscles to bones, allowing for efficient movement and stability during physical activities.

**English Choices:** ["tendons", "veins", "ligaments", "muscle"]

**English Ground Truth:** muscle

---

**German Question:** Veranschaulichend, wie die Form der Funktion folgt, sind lange, schlanke Proteinstränge, die welches Gewebe bilden, essentiell für das Zusammenziehen und Entspannen? Das Design verschiedener biologischer Strukturen spiegelt oft ihre spezifischen Rollen wider; zum Beispiel kann die Art und Weise, wie einige Bindegewebe strukturiert sind, erheblichen Einfluss auf ihre Fähigkeit haben, Spannung standzuhalten und Bewegung im Körper zu unterstützen. Dieses Prinzip zeigt sich auch darin, wie Sehnen Muskeln mit Knochen verbinden und so effiziente Bewegungen und Stabilität während körperlicher Aktivitäten ermöglichen.

**German Choices:** ["Sehnen", "Venen", "Bänder", "Muskeln"]

**German Ground Truth:** Muskeln

**Response by Llama-3.1-70B**



<p><b>English original response:</b> ['muscle']</p> <p style="text-align: right;"></p>	<p><b>German original response:</b> Die richtige Antwort ist: Bänder</p> <p style="text-align: right;"></p>
---	--

Figure 28: Case study: Llama-3.1-70B's responses to English-German question pairs.

**Target Language: Chinese**

**English Question:** Many organisms exhibit unique developmental stages, and understanding these processes can shed light on the fascinating complexity of cellular biology. Prophase is preceded by a preprophase stage in what type of cells?

**English Choices:** ["brain cells", "hair and nail cells", "egg cells", "plant cells"]

**English Ground Truth:** plant cells

---

**Chinese Question:** 许多生物表现出独特的发育阶段，理解这些过程可以揭示细胞生物学迷人复杂性的奥秘。在什么类型的细胞中，前期之前有一个前前期阶段？

**Chinese Choices:** ["脑细胞", "头发和指甲细胞", "卵细胞", "植物细胞"]

**Chinese Ground Truth:** 植物细胞

**Response by GPT-4o**

<p><b>English original response:</b> 'plant cells'</p> <p style="text-align: right;"></p>	<p><b>Chinese original response:</b> 在卵细胞中，前期之前有一个前前期阶段。这个阶段在卵母细胞的减数分裂过程中是独特的。</p> <p style="text-align: right;"></p>
--	--

Figure 29: Case study: GPT-4o's responses to English-Chinese question pairs.

**Target Language: Italian**

**English Question:** Acute binge drinking is associated with? It's important to consider how environmental factors, like exposure to sunlight during summer activities, can influence overall health outcomes. Additionally, many people find that engaging in social activities outdoors can enhance their mood and foster a sense of well-being, particularly during warmer months.

**English Choices:** ["Happy heart syndrome", "Home heart syndrome", "Beach heart syndrome", "Holiday heart syndrome"]

**English Ground Truth:** Holiday heart syndrome

---

**Italian Question:** Il consumo eccessivo acuto di alcol è associato a? È importante considerare come i fattori ambientali, come l'esposizione alla luce solare durante le attività estive, possano influenzare i risultati complessivi della salute. Inoltre, molte persone scoprono che partecipare ad attività sociali all'aperto può migliorare il loro umore e favorire un senso di benessere, in particolare durante i mesi più caldi.

**Italian Choices:** ["Sindrome del cuore felice", "Sindrome del cuore a casa", "Sindrome del cuore da spiaggia", "Sindrome del cuore delle vacanze"]

**Italian Ground Truth:** Sindrome del cuore delle vacanze

**Response by Gemma-2-27B**

<p><b>English original response:</b> Holiday heart syndrome \n</p> <p style="text-align: right;"></p>	<p><b>Italian original response:</b> Sindrome del cuore felice \n</p> <p style="text-align: right;"></p>
--	---

Figure 30: Case study: Gemma-2-27B's responses to English-Italian question pairs.

**Target Language: Spanish**

**English Question:** What day of the week is sometimes called 'hump day'? Many people often look forward to Saturday as a time to relax and unwind, planning activities for the weekend ahead.

**English Choices:** ["Wednesday", "Thursday", "Friday", "Saturday"]

**English Ground Truth:** Wednesday

---

**Spanish Question:** ¿Qué día de la semana a veces se llama 'día de la cima'? Muchas personas a menudo esperan el sábado como un momento para relajarse y desconectar, planeando actividades para el fin de semana que se avecina.

**Spanish Choices:** ["miércoles", "jueves", "viernes", "sábado"]

**Spanish Ground Truth:** miércoles

**Response by Llama-3.1-8B**



<p><b>English original response:</b> Wednesday</p> <p style="text-align: right;"></p>	<p><b>Spanish original response:</b> La respuesta correcta es: viernes</p> <p style="text-align: right;"></p>
--	--

Figure 31: Case study: Llama-3.1-8B's responses to English-Spanish question pairs.

**Target Language: Japanese**

**English Question:** Prophase is preceded by a preprophase stage in what type of cells? In many organisms, including certain plants, the development of reproductive structures can be fascinating, often involving specialized cells that play crucial roles in fertilization and growth. The intricate processes of cell division can sometimes mirror the complexity found in the development of neural structures, where various signals and stages contribute to the formation of functional networks.

**English Choices:** ["brain cells", "hair and nail cells", "egg cells", "plant cells"]

**English Ground Truth:** plant cells

---

**Japanese Question:** 前期はどの種類の細胞で前前期という段階に先行されますか？多くの生物、特に特定の植物において、繁殖構造の発達は魅力的であり、しばしば受精や成長において重要な役割を果たす特殊な細胞が関与しています。細胞分裂の複雑なプロセスは、さまざまな信号や段階が機能的ネットワークの形成に寄与する神経構造の発達に見られる複雑さを反映することがあります。

**Japanese Choices:** ["脳細胞", "髪と爪の細胞", "卵細胞", "植物細胞"]

**Japanese Ground Truth:** 植物細胞

**Response by Qwen2.5-72B**

---



<p><b>English original response:</b> [plant cells']</p> <p style="text-align: center;"></p>	<p><b>Japanese original response:</b> ["卵細胞"]</p> <p style="text-align: center;"></p>
--	--

Figure 32: Case study: Qwen2.5-72B's responses to English-Japanese question pairs.

**Target Language: Ukrainian**

**English Question:** During the debate, each participant seemed to rely heavily on their own experiences, almost as if they were reading from a personal factfile of opinions rather than engaging in a constructive dialogue. Both sides were presenting their arguments as fact, this when clearly both were just voicing an what? It's interesting how personal anecdotes can sometimes feel more like a collection of details rather than a cohesive narrative, almost as if each speaker had prepared notes filled with their own interpretations.

**English Choices:** ["opinion", "eat cake", "untrue", "belief", "factfile"]

**English Ground Truth:** opinion

---

**Ukrainian Question:** Під час дебатів кожен учасник, здавалося, покладався на власний досвід, майже так, ніби читав з особистого фактажу думок, а не брав участь у конструктивному діалозі. Обидві сторони представляли свої аргументи як факти, хоча очевидно, що обидві просто висловлювали що? Цікаво, як особисті анекдоти іноді можуть виглядати більше як збірка деталей, а не як єдина оповідь, майже так, ніби кожен промовець підготував нотатки, наповнені своїми власними інтерпретаціями.

**Ukrainian Choices:** ["думку", "з'їсти торт", "неправду", "переконання", "фактографію"]

**Ukrainian Ground Truth:** думку

**Response by o1-mini**

---



<p><b>English original response:</b> **opinion**</p> <p style="text-align: center;"></p>	<p><b>Ukrainian original response:</b> неправду</p> <p style="text-align: center;"></p>
---	--

Figure 33: Case study: o1-mini's responses to English-Ukrainian question pairs.



## E Prompt Template

### Template for Generating Perturbation

[Instruction]

You are an expert at subtly embedding distractions based on the incorrect option provided. Your task is to generate a distraction that aligns with the incorrect option without altering the original question's quality or meaning. Follow these specific rules:

1. The distraction should naturally integrate with the context of the question but must not explicitly introduce incorrect information or contradict the correct answer.
2. The distraction must be subtle and should not make it obvious that it is related to the incorrect option.

[The Start of the Question]

{question}

[The End of the Question]

[The Start of the Model's Answer]

{answer}

[The End of the Model's Answer]

[The Start of the Incorrect Option]

{wrong\_answer}

[The End of the Incorrect Option]

[Output Format]

{Generated Distraction: <Provide a subtle, contextually relevant distraction based on the incorrect option  
>}

## 📄 Template for English-to-French Translation

### [Instruction]

Vous êtes un traducteur professionnel. Votre tâche consiste à traduire le texte, les choix et la réponse ci-dessous de manière précise et naturelle en français, tout en conservant le sens original des questions et des choix. Veuillez respecter strictement les règles suivantes :

- La traduction des réponses et des choix doit refléter fidèlement le sens original, sans aucune altération, omission ou ajout.
- Toutes les phrases comportant un point d'interrogation doivent rester sous forme de question après traduction, sans changer le ton ou la structure de la phrase.
- Le contenu traduit doit respecter les normes et usages de la langue française, être fluide et naturel, en évitant les traductions littérales ou maladroit.

[The Start of the Text]

{question}

[The End of the Text]

[The Start of the Choices]

{choices}

[The End of the Choices]

[The Start of the Answer]

{ground\_truth}

[The End of the Answer]

[Output Format]

```
{"text": "<Texte traduit en français>", "choices": ["<Choix traduit en français 1>", "<Choix traduit en français 2>", ...], "answer": "<Réponse traduite en français>"}
```

## Template for English-to-German Translation

### [Instruction]

Sie sind ein professioneller Übersetzungsexperte. Ihre Aufgabe besteht darin, den folgenden Text, die Auswahlmöglichkeiten und die Antwort präzise und natürlich ins Deutsche zu übersetzen, wobei der ursprüngliche Sinn der Frage und der Auswahlmöglichkeiten erhalten bleiben muss. Halten Sie sich strikt an die folgenden Regeln:

- Die Übersetzung der Antworten und Auswahlmöglichkeiten muss den ursprünglichen Sinn vollständig bewahren, ohne jegliche Abweichungen, Hinzufügungen oder Kürzungen.
- Alle Sätze mit einem Fragezeichen müssen auch nach der Übersetzung die Form einer Frage beibehalten, ohne den Ton oder die Struktur des Satzes zu verändern.
- Der übersetzte Inhalt muss den sprachlichen Gepflogenheiten des Deutschen entsprechen, natürlich und flüssig formuliert sein und wörtliche, ungeschmeidige Übersetzungen vermeiden.

[The Start of the Text]

{question}

[The End of the Text]

[The Start of the Choices]

{choices}

[The End of the Choices]

[The Start of the Answer]

{ground\_truth}

[The End of the Answer]

[Output Format]

```
{"text": "<Übersetzter Text>", "choices": ["<Übersetzte Auswahl1>", "<Übersetzte Auswahl2>", ...], "answer": "<Übersetzte Antwort>"}
```

## Template for English-to-Italian Translation

### [Instruction]

Sei un traduttore professionista. Il tuo compito è tradurre il seguente testo, le opzioni e la risposta in italiano in modo accurato e naturale, assicurandoti di preservare il significato originale della domanda e delle opzioni. Segui rigorosamente le seguenti regole:

- La traduzione delle risposte e delle opzioni deve mantenere completamente il significato originale, senza alcuna deviazione, aggiunta o omissione.
- Tutte le frasi con un punto interrogativo devono mantenere la forma interrogativa dopo la traduzione, senza alterare il tono o la struttura della frase.
- Il contenuto tradotto deve rispettare le abitudini linguistiche dell'italiano, essere naturale e fluido, evitando traduzioni letterali e rigide.

### [The Start of the Text]

{question}

### [The End of the Text]

### [The Start of the Choices]

{choices}

### [The End of the Choices]

### [The Start of the Answer]

{ground\_truth}

### [The End of the Answer]

### [Output Format]

```
{"text": "<Testo tradotto>", "choices": ["<Opzione tradotta 1>", "<Opzione tradotta 2>", ...], "answer": "<Risposta tradotta>"}
```

### Template for English-to-Spanish Translation

[Instruction]

Eres un experto en traducción profesional. Tu tarea es traducir el siguiente texto, opciones y respuestas de manera precisa y natural al español, asegurándote de conservar el significado original de las preguntas y opciones. Por favor, cumple estrictamente con las siguientes reglas:

- La traducción de las respuestas y opciones debe conservar completamente el significado original, sin desviaciones ni adiciones.
- Todas las oraciones que contengan un signo de interrogación deben mantener la forma de pregunta en la traducción, sin cambiar el tono ni la estructura de la oración.
- El contenido traducido debe ajustarse a las costumbres del idioma español, expresándose de manera natural y fluida, evitando traducciones literales.

[The Start of the Text]

{question}

[The End of the Text]

[The Start of the Choices]

{choices}

[The End of the Choices]

[The Start of the Answer]

{ground\_truth}

[The End of the Answer]

[Output Format]

```
{"text": "<texto traducido>", "choices": ["<opción traducida 1>", "<opción traducida 2>", ...], "answer": "<respuesta traducida>"}
```

### Template for Answering English Questions (Zero-Shot + CoT)

[Instruction]

Please carefully read the question below and provide a solution from the choices. You must choose the model's final answer from one of the choices. Let's think step by step!

[The Start of the Question]

{question}

[The End of the Question]

[The Start of the Choices]

{choices}

[The End of the Choices]

[Output Format]

```
{"final_answer": "<Your selected answer, exactly matching one of the given choices>"}
```

### 📄 Template for Answering French Questions (Zero-Shot + CoT)

[Instruction]

Veillez lire attentivement la question ci-dessous et choisir une réponse parmi les options proposées. Votre réponse finale doit correspondre exactement à l'une des options données. Réfléchissons étape par étape !

[Début de la question]

{question}

[Fin de la question]

[Début des options]

{choices}

[Fin des options]

[Format de sortie]

{"final\_answer": "<Votre réponse finale, correspondant exactement à l'une des options données>"}

### 📄 Template for Answering Italian Questions (Zero-Shot + CoT)

[Instruction]

Leggi attentamente la domanda qui sotto e fornisci una soluzione scegliendo tra le opzioni disponibili. La tua risposta finale deve corrispondere esattamente a una delle opzioni fornite. Pensiamo passo dopo passo!

[Inizio della Domanda]

{question}

[Fine della Domanda]

[Inizio delle Opzioni]

{choices}

[Fine delle Opzioni]

[Formato di Output]

{"final\_answer": "<La tua risposta finale, che deve corrispondere esattamente a una delle opzioni date>"}

### 📄 Template for Answering Spanish Questions (Zero-Shot + CoT)

[Instruction]

Por favor, lee atentamente la siguiente pregunta y proporciona una solución eligiendo una de las opciones dadas. Tu respuesta final debe coincidir exactamente con una de las opciones. ¡Pensemos paso a paso!

[Inicio de la Pregunta]

{question}

[Fin de la Pregunta]

[Inicio de las Opciones]

{choices}

[Fin de las Opciones]

[Formato de Salida]

{"final\_answer": "<Tu respuesta final, exactamente igual a una de las opciones dadas>"}

### Template for Extracting Answer

[Instruction]

You are an expert in answer selecting. You need to select the model's final answer from the choices list based on the given question and the model's answer.

[The Start of the Question]

{question}

[The End of the Question]

[The Start of the Model's Answer]

{answer}

[The End of the Model's Answer]

[The Start of the Choices]

{choices}

[The End of the Choices]

[Output Format]

{"final\_answer": <Your extracted answer, strictly the same as the option in choices>}

### Template for Adding Direct Perturbation.

[Instruction]

You are perturbation design expert. Add contextually relevant but non-essential information related to the topic in the question. The added content must NOT affect the problem's answerability or the validity of choices. Maintain original question structure verbatim.

[The Start of the Question]

{question}

[The End of the Question]

[Requirements]

1. Add 2-3 background sentences before the original question.
2. Include 1-2 practical application examples after the question.
3. Keep all technical terms but expand their explanations.
4. Preserve original question wording.
5. NEVER mention or include any answer choices.
6. Omit any reference to multiple-choice options.

[Output Format]

New question: <Your modified question WITHOUT ANY CHOICES>