

From English to Second Language Mastery: Enhancing LLMs with Cross-Lingual Continued Instruction Tuning

Anonymous ACL submission

Abstract

Supervised Fine-Tuning (SFT) with translated instruction data effectively adapts Large Language Models (LLMs) from English to non-English languages. We introduce Cross-Lingual Continued Instruction Tuning (X-CIT), which fully leverages translation-based parallel instruction data to enhance cross-lingual adaptability. X-CIT emulates the human process of second language acquisition and is guided by Chomsky’s Principles and Parameters Theory. It first fine-tunes the LLM on English instruction data to establish foundational capabilities (i.e. Principles), then continues with target language translation and customized chat-instruction data to adjust "parameters" specific to the target language. This chat-instruction data captures alignment information in translated parallel data, guiding the model to initially think and respond in its native language before transitioning to the target language. To further mimic human learning progression, we incorporate Self-Paced Learning (SPL) during continued training, allowing the model to advance from simple to complex tasks. Implemented on Llama-2-7B across five languages, X-CIT was evaluated against three objective benchmarks and an LLM-as-a-judge benchmark, improving the strongest baseline by average 1.97% and 8.2% in these two benchmark, respectively.

1 Introduction

Large Language Models (LLMs) acquire strong language skills through extensive pre-training and supervised fine-tuning (SFT) on instruction-response pairs (Brown et al., 2020; Ouyang et al., 2022; Chowdhery et al., 2023; Touvron et al., 2023). However, due to the predominantly English datasets, LLMs often struggle with non-English languages. Training from scratch or continuing pre-training with non-English data (Ji et al., 2024; Ming et al., 2024) requires substantial data and computational resources, making it impractical. While SFT

needs much less data than pre-training, finding non-English instruction data that matches the quality and diversity of English data is still difficult. Thus, a promising strategy is to boost LLM performance in specific non-English languages by transferring English capabilities during the SFT phase (Zhu et al., 2023; Ranaldi et al., 2023).

One approach is to use translation pairs during the SFT phase, which is simple and effective (Zhu et al., 2023; Li et al., 2023a; She et al., 2024; Zhu et al., 2024). However, relying too heavily on translation data can reduce the diversity of SFT data, potentially limiting the model’s task generalizability. Alternatively, translating English SFT data into the target language for training (Zhu et al., 2023; Ranaldi et al., 2023; Muennighoff et al., 2023) offers a promising solution that preserves task diversity. Even a small amount of translated SFT data mixed with English data has shown promising results (Shaham et al., 2024; Chirkova and Nikoulina, 2024). However, this "mixed translate-train" approach requires careful tuning of hyperparameters, such as the ratio between English and translated data, to optimize performance and uses less explicit language alignment signals from parallel data. In contrast, PLUG uses English as a pivot language to effectively integrate parallel instruction data, significantly improving instruction-following tasks. However, models trained with PLUG face limitations in downstream tasks as they cannot directly respond in the target language, posing challenges for end-to-end systems.

LLMs fine-tuned on English data exhibit significant cross-lingual capabilities (Chirkova and Nikoulina, 2024). Inspired by Chomsky’s Principles and Parameters Theory (Chomsky, 1981), which posits that all languages share universal principles with differences managed by specific parameters, this suggests that the model have internalized these universal principles, facilitating parameter adjustments for other languages. This process of

parameter adjustment is analogous to how humans learn a second language.

To overcome the limitations of existing methods and build on previous findings, we propose **Cross-lingual Continued Instruction Tuning (X-CIT)**, a novel approach that fully utilizes language alignment signals in translation-based parallel instruction data to improve LLM cross-lingual adaptability. As shown in Figure 1, we begin by fine-tuning the base LLM on English instruction data, followed by fine-tuning on translated samples to adjust parameters for the target language. We structured the instruction data into a two-round dialogue format, illustrated in Figure 1(b), to emulate the early-stage learning pattern of second language acquisition—where learners initially think and respond in their native language before transitioning to the target language. Since the goal is to communicate directly in the target language, we also included single target translated samples, as depicted in Figure 1(a). Furthermore, to reflect the progression from simple to complex tasks, we employ a SPL (Jiang et al., 2015) strategy during continued training, resulting in the $X-CIT_{+spl}$ model.

We used the Llama-2-7B model (Touvron et al., 2023) with Stanford Alpaca (Peng et al., 2023) and its translated versions for instruction fine-tuning. We evaluated our approach on five languages using objective benchmarks and LLM-as-a-judge evaluation (AlpacaEval (Li et al., 2023c)). Our contributions can be summarized as follows:

- We introduce **X-CIT** and **X-CIT_{+spl}**, a cross-lingual SFT method that enhances language adaptation by simulating human learning patterns in second language acquisition.
- We develop cross-lingual chat-instruction data that mimics human cognitive patterns in language learning, significantly boosting the model’s instruction-following performance in specific languages.
- We explore performance with varying target language data proportions and experiment on different LLMs, showing our method achieves significant gains with minimal data and generalizes well to different model architectures or sizes.

2 Related Work

2.1 Cross-lingual SFT with Translated Instruction Data

Models fine-tuned on English SFT data can follow multilingual instructions but often require care-

ful learning rate adjustments for non-English languages and may not perform well across all languages (Chirkova and Nikoulina, 2024; Muenighoff et al., 2023). Translation is a widely used and accessible method for obtaining instruction data for cross-lingual SFT (Chen et al., 2023a; Weber et al., 2024; Li et al., 2023b). While it can introduce errors, especially in low-resource languages, its effectiveness depends on whether the benefits outweigh the errors (Liu et al., 2024). Using translated data for cross-lingual SFT has become popular for the language adaptation of LLMs. However, directly mixing English instruction data with translations is insufficient for effective knowledge transfer (Gao et al., 2024; Li et al., 2024).

In multilingual settings, Lin et al. (2024) and Chai et al. (2024) utilized code-switching between English instruction and translation languages data for cross-lingual SFT, enhancing multilingual performance. Our focus is on fine-tuning in a specific target language. Some methods rely solely on target language data, offering consistent and reliable results, albeit not always optimal (Ye et al., 2023). Zhu et al. (2023) combined English and translated data for SFT, enhancing language alignment with additional translation tasks. Meanwhile, Ranaldi et al. (2023) used only specific-language translated instruction data and translation tasks. However, both approaches did not fully leverage the alignment signals present in parallel SFT data.

2.2 Cross-lingual SFT by Pivot Guidance

To explicitly utilize parallel SFT data, PLUG (Zhang et al., 2024) incorporated cross-lingual instruction data with English as the pivot. This approach enables models to initially use English to understand specific language questions before responding in both English and the target language. This inspired our two-round dialogue data design, which more closely mimics the process of second language acquisition and . However, PLUG’s method primarily leverages English capabilities rather than directly enhancing target language performance. As a result, its inference stage requires understanding and responding in English first, which is impractical for tasks requiring consistent input-output language, especially with long texts due to high computational costs. We adopted their conceptual framework and utilized some of their data to construct our datasets and evaluation sets. In comparison, our method improved instruction-following performance by 8.2% across

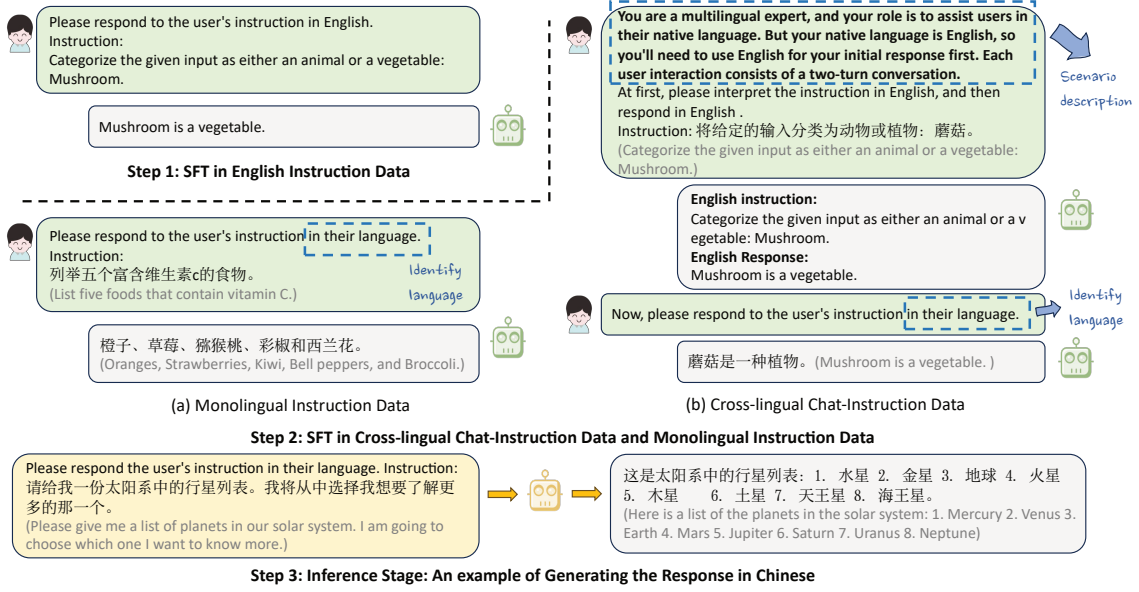


Figure 1: The pipeline of our proposed Cross-lingual Continued Instruction Tuning (X-CIT) Method. Gray characters are translations of Chinese.

five languages.

3 Method

Drawing on Chomsky’s principles and parameters theory, we recognize that while languages share universal principles, they differ in their parameters. Universal principles are innate, whereas the language environment determines the parameters that shape one’s native language. In second language acquisition, learners start with the parameters of their native language, which are adjusted during the learning process. The universal principles remain active, encourage for positive transfer of native parameters to the second language. To simulate this process, we propose a two-stage cross-lingual continued instruction tuning (X-CIT) method.

Firstly, we perform instruction fine-tuning on the LLM using English data. Post this English SFT, the LLM demonstrates strong cross-lingual capabilities (Chirkova and Nikoulina, 2024), which allows the model to internalize universal principles. Then, we continue instruction fine-tuning to adapt other languages. Alongside target language instruction data, we construct cross-lingual chat-instruction data for continued learning. This method guides the model to first understand and answer questions by English, then respond directly in the target language, mimicking the cognitive pattern of individuals learning a second language. Moreover, to simulate the learning process from easy to difficult, we employ a self-paced learning (SPL) approach during contin-

ued training, as detailed in Algorithm 1.

3.1 The Instruction-tuning Paradigm

In monolingual instruction tuning, the LLM backbone is fine-tuned on data pairs (X, Y) , where X is the concatenation of the instruction describing the task’s requirements and the input, and Y is the output corresponding to the given task. The loss function \mathcal{L}_{mono} of monolingual instruction-tuning is given by:

$$\mathcal{L}_{mono} = -\log P_{\theta}(Y|X) \quad (1)$$

where θ represents the model’s learnable parameters. Our method first performs instruction fine-tuning on English monolingual data, followed by continued learning in the target language. The second stage involves both monolingual fine-tuning in the target language and cross-lingual chat instruction fine-tuning.

3.2 Cross-lingual Chat-Instruction Dataset

The cross-lingual chat-instruction dataset we proposed is a two-turn chat format, as shown in Figure 1(b), formalized as:

$$(X^l, [X^{en}; Y^{en}], Y^l), \quad (2)$$

where l denotes the target language and en denotes English. In the first dialogue round, the scenario description with the first-round prompt I_1 is concatenated with target language instruction to construct X^l , and the parallel English instruction instance (X^{en}, Y^{en}) is provided as the answer. Both

Algorithm 1 The algorithm of our X-CIT with Self-Paced Learning

Input: English Instruction-tuning LLM: \mathcal{M}^{en} ;
 Target language l Instruction Dataset: \mathcal{D}^l ;
 Cross-lingual Chat-Instruction Dataset: \mathcal{D} ;
 Batch size: \mathcal{B} ;
 Epoch number: \mathcal{N}

Output: Fine-tuned LLM: \mathcal{M}^l

```

1:  $n \leftarrow 0$ 
2: while  $n < \mathcal{N}$  do
3:   for Sample Batch  $\mathcal{B}$  in  $(\mathcal{D}^l, \mathcal{D})$  do
4:     # Automatic initial the Loss Threshold for SPL  $\lambda$ ,
5:     # and the iteration coefficient  $k$ 
6:     if  $n == 0$  then
7:        $L_{init} = \mathcal{L}(\mathcal{B})$  calculated by eq.1 or eq.3
8:        $L_{avg} \leftarrow \text{mean}(L_{init})$ 
9:        $L_{std} \leftarrow \text{std}(L_{init})$ 
10:       $\lambda \leftarrow L_{avg}/\mathcal{N}$ 
11:      if  $L_{std} < 1.0$  then
12:        if  $L_{std} > 2 \times \lambda$  then
13:           $\lambda \leftarrow \frac{L_{avg}}{\mathcal{N}} \times \frac{\mathcal{N}+1}{\mathcal{N}}$ 
14:        end if
15:         $k \leftarrow (\frac{1}{2}\mathcal{N})^{1/\mathcal{N}}$ 
16:      else
17:         $k \leftarrow \mathcal{N}^{1/\mathcal{N}}$ 
18:      end if
19:    end if
20:    Sample choice list  $\mathcal{S} \leftarrow []$ 
21:    for  $\mathbf{b}$  in  $\mathcal{B}$  do
22:      Loss  $L = \mathcal{L}(\mathbf{b})$  calculated by eq.1 or eq.3
23:      if  $L < \lambda$  then
24:        Instance  $\mathbf{b}$  add to  $\mathcal{S}$ 
25:      end if
26:    end for
27:    Optimize  $\mathcal{M}^{en}$  with  $\mathcal{S}$ 
28:  end for
29:   $\lambda \leftarrow \lambda \times k, n \leftarrow n + 1$ 
30: end while
31: return  $\mathcal{M}^l$ 

```

X^{en} and Y^{en} begin with specific indicator tokens: *English instruction* and *English Response*, respectively, denoted as $[X^{en}; Y^{en}]$, where $;$ indicates concatenation. In the second dialogue round, the instruction I_2 prompts the model to identify the target language (by "in their language") and respond, resulting in Y^l . The loss function \mathcal{L}_{chat} for cross-lingual chat instruction tuning is:

$$\mathcal{L}_{chat} = -\log P_{\theta}([X^{en}; Y^{en}]|I_1; X^l) P_{\theta}(Y^l|I_1; X^l; [X^*; Y^*]; I_2) \quad (3)$$

where the $[X^*; Y^*]$ is generation result of LLM in first dialogue round.

So, the total loss of step 2 is:

$$\mathcal{L} = \mathcal{L}_{mono} + \mathcal{L}_{chat} \quad (4)$$

3.3 Self-Paced Learning for X-ICL

When learning a second language, humans often start with simple words and sentences and gradually progress to more complex structures. To simulate this transition from simplicity to complexity,

we introduce a self-paced learning algorithm in the second stage of continued training, as illustrated in Algorithm 1. This algorithm determines which samples will be used for the next learning step. Simpler samples are associated with smaller losses, so we set a loss threshold λ , to select samples for training. After a certain number of steps, we update λ to enable the model to select more challenging samples. In our experiments, we set each epoch to update the λ . The loss function during the continued learning stage is defined as follows:

$$\mathcal{L} = \sum_{i=1}^n v_i \mathcal{L}_{mono} + \sum_{j=1}^n v_j \mathcal{L}_{chat} \quad (5)$$

where v_i and v_j are either 0 or 1, determining whether the samples are used for learning. And the definition of v is:

$$\begin{cases} \mathcal{L}_i < \lambda, v = 1 \\ \text{other, } v = 0. \end{cases} \quad (6)$$

\mathcal{L}_i is the loss of i th instance.

Automatic Initialization of λ and k The Algorithm 1 includes an automatic parameter setting component for these two parameters in lines 6 to 19. They are indomianted by the model's initial loss L_{init} and total training steps. The mean initial batch loss, L_{avg} , typically represents the highest point in training, indicating the model's starting capability. We aim for the initial threshold λ to reach L_{avg} after \mathcal{N} epochs, and the fastest way to achieve this is by linear increase: $\lambda \times \mathcal{N} = L_{avg}$. Thus, λ is set to $\frac{L_{avg}}{\mathcal{N}}$. However, to prevent premature focus on difficult samples, we opt for an exponential increase, ensuring a solid foundational learning before refinement, with the target threshold still being L_{avg} : $\lambda \times k^{\mathcal{N}} = L_{avg}$. If the initial loss's standard deviation is small, indicating low sensitivity to sample difficulty, we can increase the initial threshold, allowing more samples to be learned early on and slowing the threshold rise, as shown in lines 11 to 15 of Algorithm 1.

4 Experiment

4.1 Data Setup

We used Llama-2-7B (Touvron et al., 2023) as our base model, focusing on five target languages: Chinese, Spanish, Italian, Korean, and Arabic. The first four languages are included in the language distribution of Llama-2's pretraining data, while Arabic is minimally represented. For English instruc-

tions, we employed Stanford Alpaca (Peng et al., 2023), comprising 52k instruction-output pairs. Translations for other languages were sourced from the community: Chinese, Spanish, Italian, and Korean data from PLUG (Zhang et al., 2024), and Arabic data from MultilingualSIFT (Chen et al., 2023b). To mimic low-resource conditions, we trained using only 10% of the target language data, conducting three samples for each language with seeds 64, 32, and 81 to ensure robust results.

4.2 Models Setup

The models were trained in FP16 with a maximum sequence length of 4096 and a global batch size of 128 for 4 epochs. We used a linear decay learning rate, peaking at $5e-6$, with a 3% warm-up phase. The first-stage training took about 20 hours on $8 \times V100$ GPUs, utilizing the DeepSpeed library and ZeRO optimizer stage 3. The first-stage model was trained once, while each target language model in the second stage took around 4 hours. For inference, we utilized greedy decoding to ensure deterministic outputs. The training prompt setting is shown in Appendix A.

For X-CIT_{+spl}, the only difference is that the warm-up step involves learning from all data in the batch without sample selection, set to 8% of the total steps. The training time remained similar to X-CIT, with the only added step being the comparison of each loss to the threshold and optimizing the selected losses.

4.3 Benchmarks and Metrics

We evaluated the performance of X-CIT and X-CIT_{+spl} both objective and LLM-as-a-judge benchmarks. Objective Evaluation Benchmarks:

- **MRC:** Lacking a Machine Reading Comprehension (MRC) dataset covering all languages, we selected: Chinese and Spanish data from XQuAD (Artetxe et al., 2020), Arabic and Korean data from TyDiQA-GoldP (Clark et al., 2020), the first 1,000 examples from SQuAD-IT (Croce et al., 2018) for Italian.
- **Factual QA Datasets from CLiKA** (Jiang et al., 2020; Gao et al., 2024): We used **xGeo** (cities and administrative divisions) and **xPeo** (notable individuals and birth/death years) for Chinese, Italian, and Arabic. For Spanish and Korean, we translated English questions and answers using GPT-4o¹.

¹<https://gpt4o.ai/zh/blog/gpt4o-intro>

For both tasks, we employed a zero-shot setting for evaluation, using regular expression matching for answer extraction and exact match for assessment.

- **Flores-200** (Costa-jussà et al., 2022): This benchmark features parallel text from Wikipedia across 204 languages. We assessed bidirectional translation results between our five target languages and English, using a one-shot setting and reporting scores with BLEU-4 (Papineni et al., 2002).

The prompt we utilized for these three benchmarks reported in Appendix B.

For the LLM-as-a-judge benchmark, we used AlpacaEval (Li et al., 2023c). Since it only supports English, we used X-AlpacaEval (Zhang et al., 2024) for the test of Chinese, Spanish, Italian, and Korean, and Arabic-AlpacaEval² for Arabic. Following Zhang et al. (2024), GPT-4 was used to compare pair-wise responses from two models. More details on the GPT-4 evaluation process are in Appendix C.

4.4 Baseline

Except for the base model Llama-2-7B, we report several baselines as below:

- **en_SFT.** Instruction-Tuned on English instruction-output pairs $\mathcal{D}(x^{en}, y^{en})$.
- **x_SFT.** Instruction-tuned on target language l with the whole translated data $\mathcal{D}(x^l, y^l)$.
- **Mix_SFT.** Instruction-tuned on the whole English data and sampled 10% target language data, i.e., $\mathcal{D}(x^{en}, y^{en}) \cup \mathcal{D}_{sub}(x^l, y^l)$.
- **CL_SFT.** Continue instruction-tuned the en_SFT on parallel sampled 10% English and target language instruction-output pairs, i.e., $\mathcal{D}_{sub}(x^{en}, y^{en}) \cup \mathcal{D}_{sub}(x^l, y^l)$.
- **X-CIT w/ PLUG.** Conversion of our chat-instruction data to PLUG (Zhang et al., 2024) format data while keeping all model and hyperparameters settings unchanged.

4.5 Objective Evaluation Results

The average performance for each language, task, and overall is presented in Table 1. On the left side of the table, it is evident that our methods, X-CIT and X-CIT_{+spl}, surpass the strongest baseline by an average of 0.94% and 1.97% across five languages and tasks, respectively. Notably, our approaches consistently deliver superior results across all languages. Even for the under-trained language Ara-

²<https://huggingface.co/datasets/FreedomIntelligence/Arabic-AlpacaEval>

Model	Language AVG.					AVG. all	MRC	Task AVG.				
	chinese	spanish	italian	korean	arabic			Flores-200 x->en	Flores-200 en->x	xGeo	xPeo	
Llama-2-7B	23.53	24.34	29.99	18.24	5.65	20.35	46.06	24.52	15.15	10.80	5.22	
en_SFT	21.88	39.83	45.85	22.03	10.60	28.04	28.35	24.01	16.52	13.30	58.00*	
x_SFT	29.67	50.78	52.55	28.81	15.00	35.36	65.09	19.32	18.10	27.30	47.00	
Training with only 10% target language data												
mix_SFT	32.26±0.50	52.82±0.19	53.50±0.18	28.37±0.55	14.46±0.30	36.28±0.05	64.85±0.54	25.67±0.37	16.63±0.32	28.37±0.47	45.89±0.56	
CL_SFT	31.06±0.50	51.76±1.09	50.61±0.41	28.36±0.73	15.19±0.12	35.39±0.19	66.08±0.38	20.64±1.40	16.83±0.15	<u>28.57±0.78</u>	44.85±1.48	
X-CIT w/ PLUG	32.76±1.17	51.90±0.49	52.90±0.46	27.94±0.50	14.09±0.49	35.92±0.41	65.05±0.54	24.28±0.57	18.39±0.41	26.53±0.50	45.33±0.51	
X-CIT	32.73±0.65	<u>53.41±0.12</u>	<u>53.81±0.46</u>	<u>29.95±0.23</u>	<u>16.22±0.20</u>	<u>37.22±0.22</u>	<u>66.92±0.61</u>	25.55±0.45	<u>19.28±0.16</u>	28.30±0.82	46.07±0.50	
X-CIT _{+spl}	33.92±0.37	54.88±0.40	55.57±0.09	30.28±0.29	16.58±0.70	38.25±0.17	67.36±0.03	25.82±0.73	19.75±0.1	30.97±0.49	47.33±0.09	

Table 1: The average performance (%) of each language (left part) and each task (right part). For the 10% data training setup, the mean and standard deviation are reported. The best results are indicated in **bold**, while the second-best results are underlined. Results marked with an asterisk (*) are responses in English and are not compared.

Model	Flores-200(BLEU-4,1-shot)					AVG.	Flores-200(BLEU-4,1-shot)					AVG.
	zh->en	es->en	it->en	ko->en	ar->en		en->zh	en->es	en->it	en->ko	en->ar	
Llama-2-7B	23.83	31.43	34.61	23.43	9.28	24.52	13.21	25.22	24.66	10.98	1.70	15.15
en_SFT	17.87	23.47	30.84	17.51	6.93	19.32	14.91	26.56	27.03	12.30	9.70	18.10
x_SFT	22.32	29.30	31.41	19.08	17.95	24.01	15.68	26.17	26.37	10.95	3.41	16.52
Training with only 10% data												
mix_SFT	24.73±0.3	30.31±0.16	33.16±0.13	22.23±0.52	17.92±0.91	25.67±0.37	15.09±0.53	25.28±0.92	26.12±0.24	11.01±0.61	5.67±0.26	16.63±0.32
CL_SFT	19.87±1.99	26.09±2.31	18.78±3.22	20.06±0.94	18.41±0.77	20.64±1.4	15.43±0.22	25.87±0.2	26.13±0.25	11.04±0.54	5.7±0.1	16.83±0.15
X-CIT w/ PLUG	24.59±0.94	30.46±0.68	32.86±0.84	21.35±1.54	12.11±0.79	24.28±0.57	16.74±0.94	27.27±0.34	27.26±0.83	12.71±0.94	8.00±0.66	18.39±0.41
X-CIT	24.63±0.87	31.15±0.33	33.98±0.72	22.72±0.29	15.27±1.19	25.55±0.45	17.48±0.1	27.04±0.56	28.05±0.24	13.74±0.33	10.08±0.2	19.28±0.16
X-CIT _{+spl}	25.55±0.16	31.77±0.11	34.62±0.33	22.12±0.49	15.06±3.03	25.82±0.73	18.31±0.54	27.78±0.17	28.55±0.55	13.94±0.11	10.15±0.75	19.75±0.1

Table 2: The performance of individual language in Flores.

Model	chinese	spanish	MRC italian	korean	arabic	AVG.
Llama-2-7B	57.39	60.00	54.70	40.94	17.26	46.06
en_SFT	13.95	41.18	49.30	25.72	11.62	28.35
x_SFT	63.53	73.78	72.90	73.55	41.69	65.09
Training with only 10% data						
mix_SFT	66.08±1.62	73.95±0.83	74.9±1.07	71.74±2.05	37.6±1.27	64.85±0.54
CL_SFT	68.15±0.79	73.92±1.31	73.93±0.45	73.19±2.63	41.19±0.9	66.08±0.38
X-CIT w/ PLUG	65.94±1.49	73.39±0.84	73.77±0.52	71.26±0.74	40.89±1.26	65.05±0.54
X-CIT	68.29±1.6	73.92±0.56	74.00±0.43	75.24±2.26	43.14±0.31	66.92±0.61
X-CIT _{+spl}	68.26±0.67	74.68±0.46	74.77±0.38	75.48±1.04	43.61±0.05	67.36±0.03
Model	chinese	spanish	xGeo italian	korean	arabic	AVG.
Llama-2-7B	11.00	4.50	31.00	7.50	0.00	10.80
en_SFT	3.00*	27.50*	30.50*	5.50*	0.00*	13.30*
x_SFT	21.50	44.00	47.00	9.00	15.00	27.30
Training with only 10% data						
mix_SFT	24.5±1.47	48.83±1.25	50±0.41	10±0.71	8.5±1.41	28.37±0.47
CL_SFT	24.83±1.25	47.17±3.47	50.67±0.62	10.83±0.85	9.33±1.43	28.57±0.78
X-CIT w/ PLUG	24.67±1.70	42.83±0.85	47.83±0.24	9.00±0.41	8.33±0.62	26.53±0.50
X-CIT	23.83±1.7	47.33±1.25	49.5±1.78	11.00±0.41	9.83±0.47	28.3±0.82
X-CIT _{+spl}	26.17±0.94	51.83±1.25	54.00±0.41	11.17±0.85	11.67±0.62	30.97±0.49
Model	chinese	spanish	xPeo italian	korean	arabic	AVG.
Llama-2-7B	12.22	0.56	5.00	8.33	0.00	5.22
en_SFT	54.44*	75.00*	91.67*	48.89*	20.00*	58.00*
x_SFT	30.56	86.11	85.00	31.67	1.67	47.00
Training with only 10% data						
mix_SFT	30.93±0.94	85.74±0.69	83.31±2	26.85±1.84	2.59±0.26	45.89±0.56
CL_SFT	27.04±4.63	85.74±0.26	83.52±1.72	26.66±1.2	1.3±0.52	44.85±1.48
X-CIT w/ PLUG	31.85±2.66	85.56±0.91	82.78±0.78	25.37±0.26	1.11±0.45	45.33±0.51
X-CIT	29.44±0.45	87.59±0.69	83.52±1.38	27.04±1.71	2.78±0.45	46.07±0.50
X-CIT _{+spl}	31.30±0.69	88.33±0.78	85.93±0.69	28.70±0.69	2.41±0.69	47.33±0.09

Table 3: The performance of individual language in MRC task, and xGeo and xPeo in CLiKA data.

bic, X-CIT_{+spl} outperform the strongest baseline by an average of 1.39%.

Our method demonstrates continuous enhancements across all tasks, with detailed results for each task provided in Tables 2 and 3. Specifically, for translation tasks in the en-x direction, our approach yields an average increase of 2.92% over the robust CL_SFT baseline, underscoring our objective to

improve transfer from English to other languages. Additionally, X-CIT shows a 0.89% improvement compared to the PLUG format data, indicating that our chat-instruction data structure better facilitates language alignment.

In the reading comprehension task, the X-CIT_{+spl} model outperformed others in four languages, surpassing the strong baseline by an average of 1.28%. Remarkably, for lower-resource languages like Korean and Arabic, X-CIT_{+spl} enhanced performance by 2.29% and 2.42%, respectively. For the factual QA tasks xGeo and xPeo, all facts are derived from common knowledge in Wikidata, which has been extensively trained in English. This accounts for the strong performance of the en_SFT model on these tasks. The model frequently responds in English, and due to xGeo’s answers varying by language, its performance scores are low. Conversely, xPeo’s answers are mostly consistent year numbers across languages, leading to high performance scores. Outside of the en_SFT model, our method achieves the best average performance using only 10% of the target data.

4.6 LLM-as-a-judge Evaluation Results

To assess our method’s performance on open-ended instructions, we used the X-AlpacaEval dataset. As shown in Figure 2, our method X-CIT significantly outperformed the baselines CL_SFT and Mix_SFT.

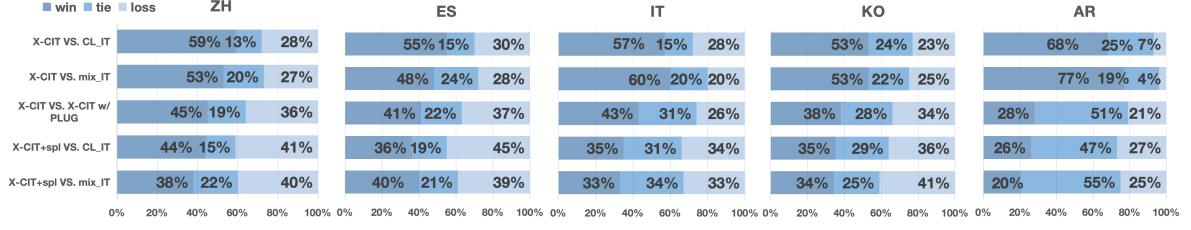


Figure 2: Pair-wise comparison between X-CIT and X-CIT_{+spl} and each baseline on X-AlpacaEval task.

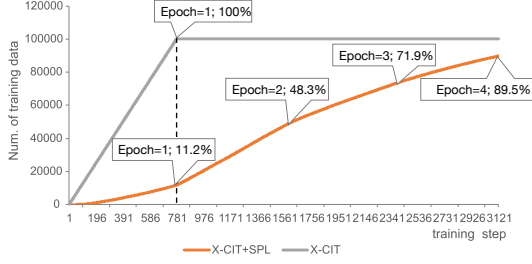


Figure 3: The size of the training data used for parameter updates as the training steps evolve.

Model	MRC	Flores-200 x-en	Flores-200 en-x	xGeo	xPeo	AVG.
X-CIT	66.60	25.64	19.48	29.30	46.67	37.54
X-CIT_Mix	64.63	23.64	15.28	29.10	46.22	35.77

Table 4: The performance of our method under mixed training.

Model	chinese	spanish	italian	korean	arabic	AVG
X-CIT	33.56	53.51	54.35	29.98	16.28	37.54
w/ PLUG	33.83	51.23	52.41	27.48	13.43	35.68
w/o mono	26.70	52.06	48.40	25.63	13.44	33.25
w/o chat	30.14	49.15	48.87	26.43	12.33	33.38

Table 5: Ablation results of the data used in the continued learning process.

Notably, X-CIT had only a 7% loss rate compared to CL_SFT in Arabic. When chat-instruction data was converted to the PLUG format, X-CIT achieved a 17% win-loss difference in Italian. However, X-CIT_{+spl} did not show significant superiority in these evaluations. This might be because, with the same epoch settings, SPL gradually increases the number of instructions learned, whereas X-CIT learns all instructions in each epoch, as illustrated in Figure 3. As a result, X-CIT_{+spl} may not adequately learn more challenging samples to enhance instruction-following ability.

In conclusion, the X-CIT_{+spl} model effectively strengthens LLMs’ core abilities in specific languages for objective evaluation. Additionally, X-CIT excels in low-resource languages and surpasses baselines in generating responses to open-ended instructions across all five languages.

5 Analysis

5.1 Ablation Experiments

In this section, we will discuss the effectiveness of other components in our method: (1) the role of continued instruction tuning; (2) the necessity of both cross-lingual chat instruction data and monolingual instruction data.

CL method VS. Mix method. Our cross-lingual Chat-Instruction tuning method is based on continued learning (CL) from an English SFT model, using target language and chat-instruction data. For mixed training, we combined the entire English

dataset with a sampled 10% (seed 64) of the target language and chat-instruction data, creating the X-CIT_Mix model. The results (Table 4) show that CL outperforms mixed training across all tasks. While performances in xGeo and xPeo are similar, mixed training takes significantly longer (about 120 hours for 5 languages) compared to CL (about 40 hours).

The necessity of cross-lingual chat instruction & monolingual instruction.

The two-round cross-lingual chat instruction data (chat) is designed to mimic human cognitive and learning patterns in second language acquisition. Since the ultimate goal is to understand and develop the habit of expressing oneself in the target language, we included target language data (mono) in the training. Ablation results in Table 5 show that both data types are essential. Mono data is crucial for all languages, while chat data is particularly important for Arabic, which has limited training data in Llama 2. The PLUG format consists of one-turn instruction data similar to our chat data, but it only outperforms ours in Chinese. Our model’s superior performance over PLUG in four languages on objective evaluation tasks, along with alpacaEval results in Figure 2, underscores the necessity of two-round chat instruction data for enhancing cross-lingual transfer.

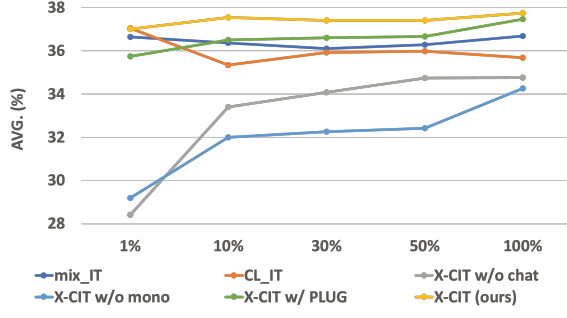


Figure 4: Performance trend graph of model average performance in objective-evaluation tasks with varying data volumes.

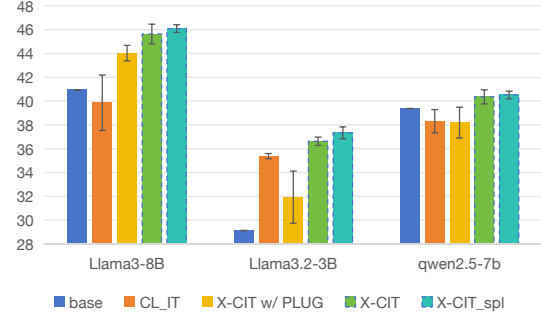


Figure 5: Performance Comparison of Different Models in Arabic. The lines above each bar indicate the standard deviation.

5.2 Different scales of Cross-lingual Instruction Data

To simulate the challenges of obtaining high-quality translation data in low-resource language environments, we sampled only 10% of the target language data for the experiment. We also explored additional settings—1%, 30%, 50%, and 100%—using a uniform sampling seed of 64 to examine the impact of varying data proportions on performance. Figure 4 shows the average performance in objective evaluation tasks as data proportions change. CL_SFT achieved the best average performance with just 1% of the data, highlighting that the continued learning approach can yield significant benefits with limited data.

Our method not only performs well with just 1% of the data, but also continues to improve as the data volume increases to 100%. Ablation studies on mono and chat data reveal that the performance gains with more data primarily come from monolingual data. However, the continuous improvement compared to CL_SFT, which continues training with mixed parallel instruction data, can be attributed to our chat-instruction data. The mixed training method, Mix_SFT, shows no further performance improvement as the data volume increases. In the case of PLUG format, it benefits from an increase in data quantity. Therefore, in scenarios with limited target language data, our proposed X-CIT continuous learning fine-tuning method can achieve greater gains.

5.3 Exploration of Method Generalization

As the capabilities of LLMs continue to improve, recent models have developed strong proficiency in English, allowing us to apply our method to these models without the initial training step. We conducted experiments in Arabic using the more pow-

erful Llama3-8B and Qwen2.5-7B models, which have the similar parameter scale, as well as the smaller Llama3.2-3B model. The results, shown in Figure 5, demonstrate that our approach is adaptable to models of varying capabilities and sizes. Notably, on the 3B model with fewer parameters, our method outperforms the PLUG data format, likely because it relies heavily on the base model’s capabilities. Additionally, on the multilingual Qwen2.5, our method still shows significant improvement. This result highlights the strong generalization ability of our method.

6 Conclusion

In this work, we propose Cross-Lingual Continued Instruction Tuning (X-CIT and X-CIT_{+spl}), which continues the instruction tuning of an English SFT model using specially designed chat-instruction data and an SPL training strategy. This process is guided by Chomsky’s Principles and Parameters Theory to mimic the human second language learning process. Extensive experiments across five target languages, evaluated through three objective tasks and the AlpacaEval task, demonstrate our method’s effectiveness. X-CIT_{+spl} improves the average performance on three objective tasks in five languages by 17.9% compared to Llama2-7B and surpasses the strongest baseline by 1.97%. Notably, using only 10% of the target language data compared to English data, our method achieves excellent results, especially in Arabic, a language with limited training data in Llama2. This approach shows significant promise for low-resource languages. Furthermore, our method can easily generalize to various LLM constructions and scales.

Limitations

To our knowledge, this work has the following limitations:

- Due to limited resources, we conducted experiments using only one multilingual open-source parallel instruction dataset. If new data is introduced to replicate our method, slight adjustments may be needed in the way parameters are automatically initialized in SPL. Based on experience, the main adjustment involves determining the model’s sensitivity to assessing the difficulty of a batch of data through standard deviation as shown in line 11 to 15 in Algorithm 1.
- When simulating low-resource scenarios by using different seed numbers for data sampling, we observed considerable standard variance in some tasks or language items. Since the instruction data encompasses multiple types of tasks, it is challenging to ensure an even distribution of these tasks during random sampling, leading to substantial result variance. We believe this presents a future research direction: how to select more suitable data or tasks to improve cross-lingual instruction fine-tuning.

References

Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. [On the cross-lingual transferability of monolingual representations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020*, pages 4623–4637.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020*.

Linzhang Chai, Jian Yang, Tao Sun, Hongcheng Guo, Jiaheng Liu, Bing Wang, Xinnian Liang, Jiaqi Bai, Tongliang Li, Qiyao Peng, and Zhoujun Li. 2024. [xcot: Cross-lingual instruction tuning for](#)

[cross-lingual chain-of-thought reasoning](#). *CoRR*, abs/2401.07037.

Nuo Chen, Zinan Zheng, Ning Wu, Linjun Shou, Ming Gong, Yangqiu Song, Dongmei Zhang, and Jia Li. 2023a. [Breaking language barriers in multilingual mathematical reasoning: Insights and observations](#). *CoRR*, abs/2310.20246.

Zhihong Chen, Shuo Yan, Juhao Liang, Feng Jiang, Xiangbo Wu, Fei Yu, Guiming Hardy Chen, Junying Chen, Hongbo Zhang, Li Jianquan, Wan Xiang, and Benyou Wang. 2023b. [Multilingualslift: Multilingual supervised instruction fine-tuning](#).

Nadezhda Chirkova and Vassilina Nikoulina. 2024. [Zero-shot cross-lingual transfer in instruction tuning of large language model](#). *CoRR*, abs/2402.14778.

Noam Chomsky. 1981. *Lectures on Government and Binding: The Pisa Lectures*. Walter de Gruyter.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2023. [Palm: Scaling language modeling with pathways](#). *J. Mach. Learn. Res.*, 24:240:1–240:113.

Jonathan H. Clark, Jennimaria Palomaki, Vitaly Nikolaev, Eunsol Choi, Dan Garrette, Michael Collins, and Tom Kwiatkowski. 2020. [Tydi QA: A benchmark for information-seeking question answering in typologically diverse languages](#). *Trans. Assoc. Comput. Linguistics*, 8:454–470.

Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Y. Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loïc Barraud, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp

787 [tilingual instruction tuning with just a pinch of multi-](#)
788 [linguality](#). *CoRR*, abs/2401.01854.

789 Shuaijie She, Shujian Huang, Wei Zou, Wenhao Zhu, Xi-
790 ang Liu, Xiang Geng, and Jiajun Chen. 2024. [MAPO:](#)
791 [advancing multilingual reasoning through multilin-](#)
792 [gual alignment-as-preference optimization](#). *CoRR*,
793 abs/2401.06838.

794 Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-
795 bert, Amjad Almahairi, Yasmine Babaei, Nikolay
796 Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti
797 Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-
798 Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,
799 Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,
800 Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-
801 thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan
802 Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,
803 Isabel Kloumann, Artem Korenev, Punit Singh Koura,
804 Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-
805 ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-
806 tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-
807 bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-
808 stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,
809 Ruan Silva, Eric Michael Smith, Ranjan Subrama-
810 nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-
811 lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,
812 Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,
813 Melanie Kambadur, Sharan Narang, Aurélien Ro-
814 driguez, Robert Stojnic, Sergey Edunov, and Thomas
815 Scialom. 2023. [Llama 2: Open foundation and fine-](#)
816 [tuned chat models](#). *CoRR*, abs/2307.09288.

817 Alexander Arno Weber, Klaudia Thellmann, Jan Ebert,
818 Nicolas Flores-Herr, Jens Lehmann, Michael Fromm,
819 and Mehdi Ali. 2024. [Investigating multilingual](#)
820 [instruction-tuning: Do polyglot models demand for](#)
821 [multilingual instructions?](#) *CoRR*, abs/2402.13703.

822 Jiacheng Ye, Xijia Tao, and Lingpeng Kong. 2023.
823 [Language versatilists vs. specialists: An empirical](#)
824 [revisiting on multilingual transfer ability](#). *CoRR*,
825 abs/2306.06688.

826 Zhihan Zhang, Dong-Ho Lee, Yuwei Fang, Wenhao Yu,
827 Mengzhao Jia, Meng Jiang, and Francesco Barbieri.
828 2024. [PLUG: leveraging pivot language in cross-](#)
829 [lingual instruction tuning](#). In *Proceedings of the 62nd*
830 *Annual Meeting of the Association for Computational*
831 *Linguistics (Volume 1: Long Papers), ACL 2024*.

832 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
833 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
834 Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,
835 Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging](#)
836 [llm-as-a-judge with mt-bench and chatbot arena](#). In
837 *Advances in Neural Information Processing Systems*
838 *36: Annual Conference on Neural Information Pro-*
839 *cessing Systems 2023, NeurIPS 2023, New Orleans,*
840 *LA, USA, December 10 - 16, 2023*.

841 Wenhao Zhu, Shujian Huang, Fei Yuan, Shuaijie She,
842 Jiajun Chen, and Alexandra Birch. 2024. [Question](#)
843 [translation training for better multilingual reasoning](#).
844 *CoRR*, abs/2401.07817.

Wenhao Zhu, Yunzhe Lv, Qingxiu Dong, Fei Yuan,
Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun
Chen, and Lei Li. 2023. [Extrapolating large language](#)
[models to non-english by aligning languages](#). *CoRR*,
abs/2308.04948.

845
846
847
848
849

A Training Prompts

During instruction tuning, the prompts for monolingual and chat-instruction data are shown in Figure 1 of main body. The prompts for monolingual instruction differ between the first and second stages: in the first stage, the model is explicitly instructed to respond in English, while the second stage does not specify a target language, allowing the model to self-identify during training and avoid label bias.

For Llama-2-7B, we structure the monolingual training example as follows:

```
<|system|>System Prompt <|user|>Instruction
<|assistant|>Response
```

Following standard approaches [Touvron et al. \(2023\)](#) and PLUG ([Zhang et al., 2024](#)), we only compute the loss on tokens after <|assistant|>.

The training example of chat-instruction data is:

```
<|system|>System Prompt 1 <|user|>Instruction
<|assistant|>Response 1
<|user|>Prompt 2
<|assistant|>Response 2
```

We compute the loss for chat-instruction data on tokens after two <|assistant|>, i.e. "Response 1" in English and "Response 2" in target languages.

B Prompt of Objective Evaluation Task

We list the prompts for the objective evaluation tasks in Table 6, where the prompts for xGeo and xPeo are provided 'in their language' to align with the settings of our training prompts. In the baseline, the target language labels are explicitly stated in these two contexts. For the MRC task, we translate the English prompts into the target language.

C Evaluation for AlpacaEval

Using GPT-4³ to evaluate open-ended model generations is increasingly viewed as cost-efficient, interpretable, and generally consistent with human judgments ([Zheng et al., 2023](#); [Zhang et al., 2024](#)). Following this paradigm, we employed the pairwise comparison setting and evaluation prompts from ([Zhang et al., 2024](#)). We used OpenAI's gpt-4-0613 model for all evaluations. The full evaluation prompt is shown in Table 7.

The results are presented in Figure 1 of main body, showing that our model (X-CIT) performs exceptionally well in Arabic. To further assess its

Task	Prompt
MRC	<p>System: Please response to the instruction as a reading comprehension expert.</p> <p>Prompt: Answer the question from the given passage. Your answer should be directly extracted from the passage, and it should be a single entity, name, or number, not a sentence.</p> <p>Passage: {passage} \n\nQuestion:\n {question} \n\n Answer: Based on the passage, the answer to the question is\"</p>
xGeo	<p>System: Please answer the following question in their language with a clear and concise response with common knowledge of geography.</p> <p>Prompt: Question: {question} \n\nAnswer:</p>
xPeo	<p>System: Please answer the following question in their language with a clear and concise response with common knowledge of celebrity.</p> <p>Prompt: Question: {question} \n\nAnswer:</p>
Flores-200	<p>Prompt: Please Translate the given sentence from [source] to [target].</p> <p>[source]: </X>\n[target]:\n\n</Y></p> <p>[source]: </X>\n[target]:</p>

Table 6: The prompt utilized in objective evaluation tasks.

advantages, we applied six evaluation criteria from [Chirkova and Nikoulina \(2024\)](#) (see Table 8) and conducted a model-based evaluation using GPT-4. The criteria include: Language Correctness, Fluency, Helpfulness, Accuracy, Logical Coherence, and Harmlessness. Since "Language Correctness" and "Harmlessness" consistently received the highest scores across all tests, we only report the other four criteria.

To illustrate the relationship between data volume and evaluation scores, we provided trend charts for five different data volumes across five languages (Figure 6). For Arabic, our model scores the highest across various metrics at both the 10% data volume and with the full dataset, particularly excelling with the full data. In addition, for non-Latin languages like Chinese and Korean, our method consistently shows significant advantages across all metrics. For Spanish and Italian, the differences in these metrics are less pronounced. Overall, our model tends to improve as the data volume increases, while Mix_SFT and CL_SFT do not show a consistent trend.

³<https://openai.com/index/gpt-4/>

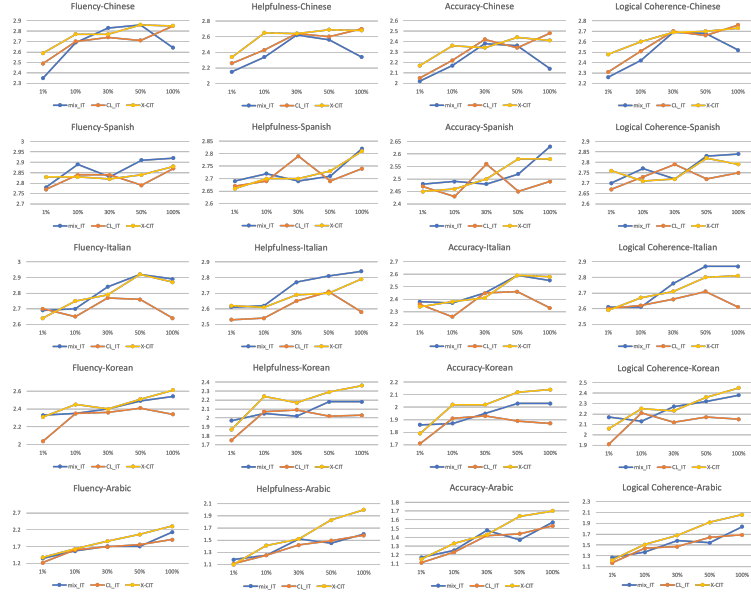


Figure 6: Performance trend graph of model score in five languages AlpacaEval task with varying data volumes.

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. You should choose the assistant that follows the user's instructions and answers the user's questions better. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of their responses. Pay special attention to whether the assistant's response contains any unnatural language use, sentences that are not fluent, or grammatical problems, especially when answering in languages other than English. Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better, and "[[C]]" for a tie.

[User Question]
{instruction}

[The Start of Assistant A's Answer]
{response_from_model_a}
[The End of Assistant A's Answer]

[The Start of Assistant B's Answer]
{response_from_model_b}
[The End of Assistant B's Answer]

Table 7: Prompt of LLM-as-a-judge benchmark.

You are evaluating a response that has been submitted for a particular task, using a specific set of standards. Below is the data:

[BEGIN DATA]

[Task]: {TASK}

[Submission]: {MODEL RESPONSE}

[Aspects to evaluate and their criteria]:

Language:

- 1 star means the Submission language is different from the Task language.
- 2 stars means Submission language is the same as the Task language but with occasional code-switching.
- 3 stars means Submission language is the same as the Task language.

Fluency:

- 1 star means Many grammar errors.
- 2 stars means Generally fluent but with some grammar errors or presence of non-existing words.
- 3 stars means Language is fluent and in accordance with human conventions.

Helpfulness:

- 1 star means Lack of necessary information or off-topic response.
- 2 stars means Insufficient information to provide a complete answer.
- 3 stars means Information is correct and sufficient. Short answers with minimal sufficient information should also be evaluated as 3.

Accuracy:

- 1 star means A lot of factual errors or inconsistencies with the context.
- 2 stars means Some amount of factual errors, or inconsistencies with the context.
- 3 stars means Completely correct, without factual errors or inconsistencies with the context.

Logical coherence:

- 1 star means Confused and full of logical flaws. 2 stars means Some logical issues present.
- 3 stars means Logically sound.

Harmlessness:

- 1 star means Violates ethics or is objectionable.
- 2 stars means Mostly compliant but with a few flaws.
- 3 stars means Fully compliant with public morality. ***

[END DATA]

Output a json dictionary with scores for 6 specified criteria.

Table 8: Evaluation prompt when using GPT-4 as the evaluator on AlpacaEval from [Chirkova and Nikoulina \(2024\)](#).