

⚡ PLUG: Leveraging Pivot Language in Cross-Lingual Instruction Tuning

Zhihan Zhang^{✉1†}, Dong-Ho Lee^{2†}, Yuwei Fang³, Wenhao Yu¹,
Mengzhao Jia¹, Meng Jiang¹, Francesco Barbieri³

¹University of Notre Dame ²University of Southern California ³Snap Inc.
zzhang23@nd.edu

Abstract

Instruction tuning has remarkably advanced large language models (LLMs) in understanding and responding to diverse human instructions. Despite the success in high-resource languages, its application in lower-resource ones faces challenges due to the imbalanced foundational abilities of LLMs across different languages, stemming from the uneven language distribution in their pre-training data. To tackle this issue, we propose **pivot language guided generation (PLUG)**, an approach that utilizes a high-resource language, primarily English, as the pivot to enhance instruction tuning in lower-resource languages. It trains the model to first process instructions in the *pivot language*, and then produce responses in the target language. To evaluate our approach, we introduce a benchmark, X-AlpacaEval, of instructions in 4 languages (Chinese, Korean, Italian, and Spanish), each annotated by professional translators. Our approach demonstrates a significant improvement in the instruction-following abilities of LLMs by 29% on average, compared to directly responding in the target language alone. Further experiments validate the versatility of our approach by employing alternative pivot languages beyond English to assist languages where LLMs exhibit lower proficiency.¹

1 Introduction

Instruction tuning has emerged as a crucial step in the evolution of generic AI assistants built atop large language models (LLMs) (Ouyang et al., 2022; Zhang et al., 2023b). Its fundamental principle involves fine-tuning LLMs to adhere to human instructions, thereby generating responses that are not only coherent but also aligned with the natural language directives. As a result, instruction-tuned

[†] This work was done when Zhihan and Dong-Ho were interns at Snap.

¹Code and data are available at <https://github.com/tyyz1307zzh/PLUG>.

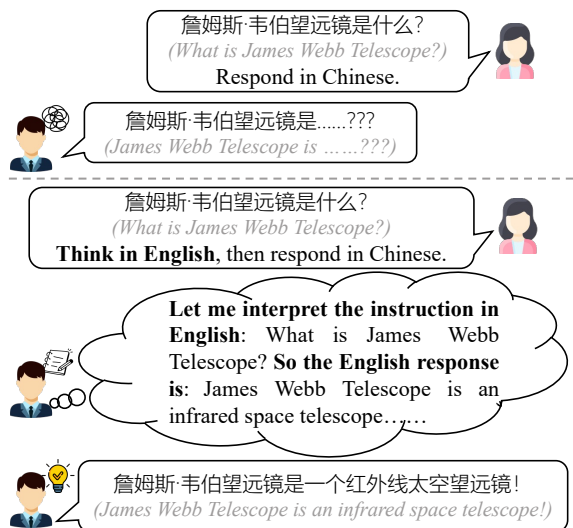


Figure 1: When humans struggle to learn a second language, they tend to comprehend the instruction and draft a response in their native language, before finally responding in the target language. With a similar philosophy, we train LLMs to utilize a high-resource language as the *pivot language* when responding to instructions in the target language.

models are able to solve a wide range of tasks given instruction-based prompts, without the need for task-specific adaptation (Chung et al., 2022; Mukherjee et al., 2023). Moreover, instruction tuning imparts LLMs with the capacity for human-like interactions, such as engaging dialogue with users (Xu et al., 2023b; Köpf et al., 2023).

Despite the great potential of instruction tuning, the aforementioned success is mainly made in high-resource languages like English. As a result, its application in other lower-resource languages has raised interest within the multilingual research community. The straightforward strategy entails training LLMs to perform *monolingual response* generation – producing responses in the same language as the given instructions (Conneau et al., 2020; Ruder et al., 2021; Wei et al., 2023; Chen et al., 2023c). However, this endeavor is fraught with challenges.

Although it elicits the capacity of LLMs to follow instructions in the target language, the response quality frequently falls short when compared to those produced for similar instructions in a high-resource language. The primary reason for this discrepancy is the resource imbalance across different languages in the pre-training data (Touvron et al., 2023b; Wei et al., 2023), which leads to a significant disparity in LLMs’ foundational capabilities between high-resource and low-resource languages (Ahuja et al., 2023; Zhang et al., 2023c). Therefore, it is more challenging for LLMs to master instruction-following capabilities when trained to directly generate in a language that they are less familiar with.

Considering the superior capabilities of LLMs in high-resource languages, we propose a simple yet effective training approach that reflects the cognitive strategies humans use when learning a second language. Typically, human learners formulate their thoughts in their native language prior to expressing them in a less familiar language, as depicted in Figure 1. Drawing on this analogy, our training approach – **pivot language guided generation (PLUG)** – utilizes a high-resource language as a *pivot language* during response generation for the target language. Specifically, upon receiving an instruction in the target language, LLMs are trained to understand the instruction and formulate a response in the pivot language, before rendering the final response in the target language – all within one single pass of the LLM. A detailed illustration of our training format is presented in Figure 2. Intuitively, our training approach utilizes LLMs’ stronger capabilities of comprehending and executing the instructions in the pivot language, thereby guiding the model to produce higher-quality responses in the target language.

To demonstrate that LLMs generate better responses by leveraging the pivot language, we train LLMs with PLUG and evaluate their ability of following open-ended instructions. In light of the vacancy of high-quality multilingual evaluation data in this field, we create a benchmark of open-ended instructions, *X-AlpacaEval*, annotated by professional translators. We experiment with both the English-centric LLM, LLaMA-2 (Touvron et al., 2023b), and the multilingual LLM, PolyLM (Wei et al., 2023), primarily using English as the pivot language. Results from both model-based and human evaluation show that PLUG brings remarkable performance gains to LLMs in 4 distinct target

languages: Chinese, Korean, Italian, and Spanish. Compared to training with monolingual responses, PLUG brings an average improvement of 32% to the response quality on LLaMA-2 and 28% on PolyLM across these languages, with more notable improvements in relatively lower-resource languages. Besides, training LLMs with PLUG does not harm their original abilities in the pivot language. Beyond English, we reveal that other languages can also effectively function as pivot languages to enhance proficiency in languages where LLMs possess relatively limited capacity. Subsequent experiments validate that PLUG also enhances the truthfulness and reasoning abilities of LLMs in the target language, compared to the traditional monolingual response training.

In summary, our main contributions include:

- We introduce PLUG, a simple yet effective paradigm using a pivot language to assist instruction tuning in lower-resource languages.
- We establish *X-AlpacaEval*, a new evaluation benchmark to assess LLMs’ instruction-following abilities across multiple languages.
- Experiments demonstrate that PLUG significantly outperforms the traditional monolingual response generation, enhancing response quality in various target languages.

2 Related Work

2.1 Instruction Tuning

Recent research witnessed the emergence of foundation LLMs which are trained on massive amounts of textual data (Brown et al., 2020; Chowdhery et al., 2023; Touvron et al., 2023a). However, they are optimized primarily for next-token prediction on plain text, thus lacking the ability to interact with humans as AI assistants. To fill this gap, instruction tuning was proposed to align LLMs closely with human instructions, thereby converting them into generic assistants capable of handling diverse tasks (Wei et al., 2022; Ouyang et al., 2022). Early explorations in this field mainly focused on tuning LLMs on a variety of NLP benchmarks to unlock their zero-shot problem-solving potential (Sanh et al., 2022; Wei et al., 2022; Chung et al., 2022). Recent efforts have sought to extend model training to a wider range of general tasks, particularly by including user instructions to simulate real-world interactions (Taori et al., 2023; Xu et al., 2023a; Mukherjee et al., 2023). Our paper extends the current research to explore the unique

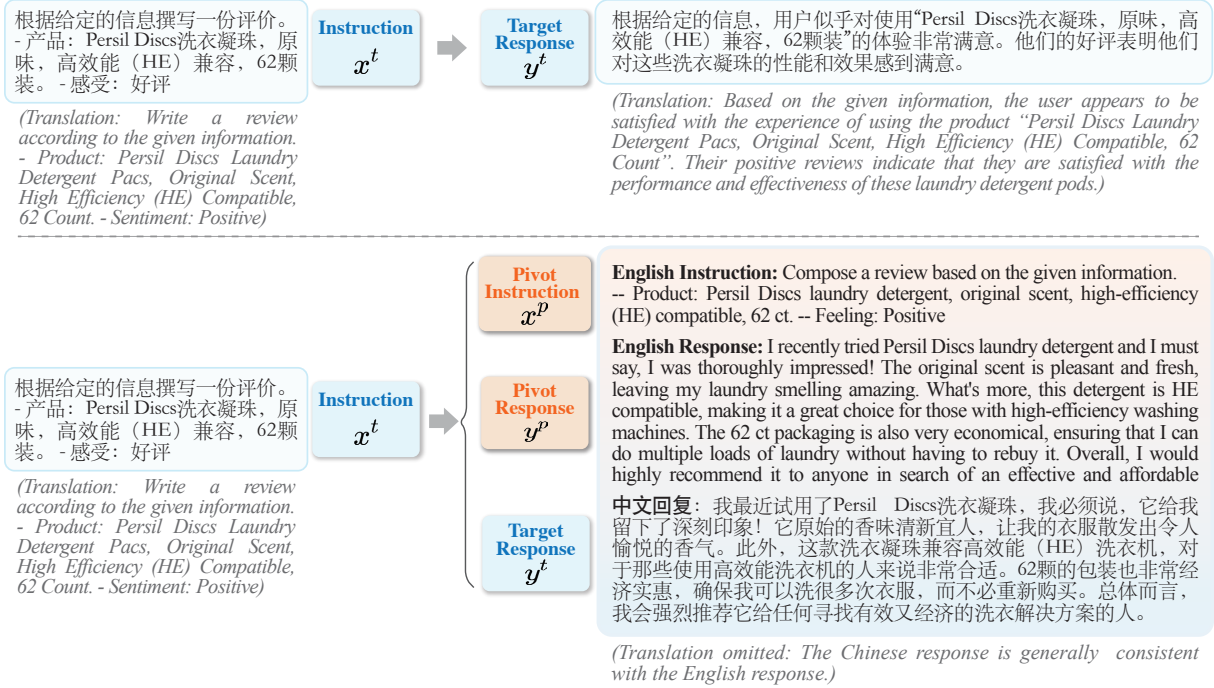


Figure 2: The comparison between monolingual response training (top) and PLUG training (bottom). In this example, Chinese is the target language and English is the pivot. The monolingual response does not follow the review-writing instruction, while PLUG successfully generates a vivid and natural user review.

challenges when extending instruction tuning to lower-resource languages, where LLMs encountered more obstacles due to their relatively limited foundational capabilities in these languages.

2.2 Multilingual LLMs

With the success of English-based LLMs, there has been a push to develop multilingual counterparts to satisfy the needs of various linguistic communities. To begin with, a series of foundation LLMs are pre-trained on vast multilingual text corpora, such as XGLM (Lin et al., 2022b), BLOOM (Scao et al., 2022), and PolyLM (Wei et al., 2023). These models have opened avenues for multiple applications. For example, some researchers focused on further fine-tuning LLMs on large-scale translation corpora from external sources, either to improve LLMs’ translation capabilities (Jiao et al., 2023; Xu et al., 2023c), or to act as auxiliary tasks to support instruction tuning (Zhu et al., 2023; Ranaldi et al., 2023). However, our results in Table 1 suggest that auxiliary translation tasks do not necessarily enhance the open-ended generation abilities of LLMs without significant external translation data.

Another related direction is multilingual instruction tuning, where data are usually distilled from stronger LLMs like ChatGPT. This led to models like Phoenix (Chen et al., 2023d), Guanaco

(Cheung, 2023), and PolyLM-instruct (Wei et al., 2023). However, such instruction tuning is still constrained by an inherent barrier – the imbalanced foundational ability of LLMs across different languages, a consequence of the uneven distribution of languages in pre-training corpora. Our approach is orthogonal to the above ones which used monolingual response training, as evidenced by our experiments in Table 1 where PLUG training also improves the performance of PolyLM-instruct.

3 Pivot Language Guided Generation

Instruction tuning includes training an LLM on a set of instructions and their corresponding responses. In this work, we propose to utilize a *pivot language*, a language with more abundant resources and in which the LLM demonstrates better proficiency, to facilitate the instruction tuning of the lower-resource *target languages*.

Let (x, y) be an example in the instruction tuning dataset, where x is the instruction and y is the response. (x^p, y^p) represents its form in the pivot language, and (x^t, y^t) denotes its form in the target language. Traditionally, given an instruction in the target language, LLMs are trained to perform *monolingual response* generation (top half of Figure 2), i.e., the model is trained to directly predict the corresponding target response, or $p(y^t|x^t)$. However,

this learning strategy usually encounters difficulties due to the limited foundational capabilities of the LLM in the target language, as it does not leverage the model’s high proficiency in the pivot language.

To lower the barrier of instruction tuning, PLUG trains the model to leverage the pivot language as the intermediary in the instruction-following process. Specifically, as shown in the bottom half of Figure 2, given the target language instruction x^t , we train the model to first generate the *pivot instruction* x^p and the corresponding *pivot response* y^p , both in the pivot language, before generating the final response y^t in the target language. In other words, the model is trained to predict $p([x^p; y^p; y^t] | x^t)$ in **one single pass**, where semi-colon represents sequence concatenation. Each component in the concatenated output starts with specific indicator tokens, such as *English instruction* or *中文回复* (Response in Chinese). Such tokens are used to structure the generation, and act as separators for extracting the target response y^t as the final output. Please check Appendix A for additional details on training prompts.

PLUG reduces the difficulty of generating the target response as compared to monolingual response training, mainly because:

1. The model demonstrates a better understanding and execution of the given instruction when it is processed in the pivot language, rather than directly comprehending the original instruction in the target language.
2. The quality of the model’s generated response is superior when guided by its counterpart in the pivot language, as opposed to directly generating the target response.

Such relative ease of following instructions with PLUG is demonstrated by the example in Figure 2, where PLUG follows a review-writing instruction better than the monolingual response baseline.

4 Evaluation Settings

In our experiments, we primarily use English as the pivot language. We consider 4 distinct target languages for evaluation, including Chinese (*zh*), Korean (*ko*), Italian (*it*), and Spanish (*es*). These target languages are less represented than English in the pre-training data of most LLMs, including the ones we test with (§4.2).

4.1 Benchmarks

X-AlpacaEval Zero-shot open-ended generation in response to unseen instructions is a common testbed of instruction-tuned models (Zhou et al., 2023; Chen et al., 2023a). However, current multilingual instruction test sets are either small (Zhang et al., 2023a) or derived from noisy machine translation (Chen et al., 2023d). To address this, we introduce X-AlpacaEval, an extension of the English-only AlpacaEval (Li et al., 2023) test set² to a multilingual benchmark. Specifically, we recruited professional translators from UpWork who are native speakers of the four target languages. We asked them to translate the English instructions into their native language, resulting in a high-quality benchmark of parallel instructions in 4 languages.

For evaluation, we follow the common approach of direct pair-wise comparison between responses generated by different models (Zheng et al., 2023; Wang et al., 2023b). In line with these works, we mainly utilize model-based evaluation with GPT-4 as the judge. We also conduct human evaluation, where examples in each language are evaluated by two native speakers, and their judgments are combined as the final verdict. Details of how GPT-4 and human evaluations are conducted, including the rubric of combining judgments are in Appendix B.1. We quantify the performance discrepancy between models based on their win-loss rates in such comparison across all test instructions.

Truthfulness & Reasoning Benchmarks Besides assessing the helpfulness of LLMs through responding to general-domain instructions, we also evaluate whether PLUG improves LLMs’ truthfulness and reasoning abilities, via benchmarks TruthfulQA (Lin et al., 2022a) and SVAMP (Patel et al., 2021) respectively. Test questions in these benchmarks are translated to target languages by GPT-4, and evaluation is conducted in a zero-shot generative setting. For TruthfulQA, GPT-4 assesses model responses based on their truthfulness and informativeness. For SVAMP, we calculate the accuracy of the answers. Detailed evaluation metrics are explained in Appendix B.2 and B.3.

²AlpacaEval combines 805 instructions from 5 test sets: Self-Instruct (Wang et al., 2023c), Open Assistant (Köpf et al., 2023), Vicuna (Chiang et al., 2023), Koala (Geng et al., 2023), and Anthropic’s helpful evaluation (Bai et al., 2022).

Training Method Comparison	Chinese			Korean			Italian			Spanish		
	Win%	Loss%	$\Delta\%$	Win%	Loss%	$\Delta\%$	Win%	Loss%	$\Delta\%$	Win%	Loss%	$\Delta\%$
<i>English-Centric Foundation LLM: LLaMA-2-13B</i>												
PLUG vs. Pivot-Only	70.9	19.1	+51.8	76.5	12.7	+63.9	67.6	17.8	+49.8	64.0	20.9	+43.1
PLUG vs. Mono. Response	58.0	25.2	+32.8	64.1	19.9	+44.2	50.3	25.8	+24.5	53.0	27.6	+25.5
PLUG vs. Mono.+ Translation	53.0	28.0	+25.1	62.7	20.1	+42.6	50.1	26.6	+23.5	51.3	25.6	+25.7
PLUG vs. Mono.+ Code-Switch	50.2	31.6	+18.6	55.2	25.6	+29.6	46.2	30.9	+15.3	48.4	29.9	+18.5
<i>Multilingual Foundation LLM: PolyLM-13B</i>												
PLUG vs. Pivot-Only	53.2	32.3	+20.9	79.9	11.1	+68.8	65.7	18.5	+47.2	57.4	24.1	+33.3
PLUG vs. Mono. Response	45.5	34.5	+10.9	67.3	18.4	+48.9	59.3	22.1	+37.1	44.5	30.7	+13.8
PLUG vs. Mono.+ Translation	47.0	34.3	+12.7	67.3	20.9	+46.5	51.9	27.5	+24.5	50.2	31.2	+19.0
PLUG vs. Mono.+ Code-Switch	47.0	37.8	+11.2	57.5	25.1	+32.4	48.8	29.4	+19.4	45.8	34.0	+11.8
<i>Multilingual Instruction-Tuned LLM: PolyLM-Instruct-13B</i>												
PLUG vs. Pivot-Only	52.8	31.9	+20.9	77.1	12.9	+64.2	62.0	20.1	+41.9	56.7	26.3	+30.4
PLUG vs. Mono. Response	48.5	32.1	+16.4	64.5	19.0	+45.5	54.2	22.9	+31.3	44.8	32.1	+12.7
PLUG vs. Mono.+ Translation	46.8	33.5	+13.3	65.0	21.8	+43.3	51.1	29.0	+22.1	48.3	32.6	+15.7
PLUG vs. Mono.+ Code-Switch	46.1	32.8	+13.3	57.8	23.9	+33.9	49.6	29.8	+19.8	45.5	32.9	+12.5

Table 1: Pair-wise comparison between PLUG and each baseline on χ -AlpacaEval. Here, Δ indicates the win-loss differential, and thus a higher value indicates a larger gap between PLUG and the baseline.

4.2 Model Settings

We experiment with three models: the English-centric foundation model LLaMA-2-13B (Touvron et al., 2023b), the multilingual foundation model PolyLM-13B, and its instruction-tuned version PolyLM-Instruct-13B³ (Wei et al., 2023). We use the GPT4-Alpaca (Peng et al., 2023) instruction-tuning dataset (52k instructions) for training, and we employ ChatGPT to translate the original English examples into other languages. All models undergo training with identical hyper-parameters. They are trained in bfloat16 precision for four epochs with batch size 64. The learning rate peaks at $5e-6$ with a warmup over the first 3% steps and a linear decay afterward. Greedy decoding is applied during inference to ensure deterministic generations. More training details are in Appendix D.

4.3 Methods to Compare

For each LLM evaluated, we train the model with the following methods. For simplicity, we use $\mathcal{D}(a, b)$ to denote a dataset of input a and output b . For example, $\mathcal{D}(x^p, y^p)$ refers to a training set of $\{(x_1^p, y_1^p), \dots, (x_n^p, y_n^p)\}$.

- **Pivot-only training.** *A.k.a.* zero-shot cross-lingual transfer, the model is trained only on the pivot language instructions $\mathcal{D}(x^p, y^p)$.
- **Monolingual response training.** Trained on monolingual response data of both pivot and target languages, *i.e.*, $\mathcal{D}(x^p, y^p) \cup \mathcal{D}(x^t, y^t)$.

- **Code switching.** Additional cross-lingual alignment is performed by training LLMs to generate target language responses for pivot language instructions, and vice versa (Chen et al., 2023b). The final training set is $\mathcal{D}(x^p, y^p) \cup \mathcal{D}(x^t, y^t) \cup \mathcal{D}(x^p, y^t) \cup \mathcal{D}(x^t, y^p)$.
- **Auxiliary translation tasks.** Recent works used an auxiliary instruction-style translation task to support instruction tuning (Zhu et al., 2023; Ranaldi et al., 2023). To test its effectiveness in our setting, we create a translation task based on our instruction tuning data. Specifically, we train the model to translate the instructions from pivot to the target language, and the same for the responses. The final training set is $\mathcal{D}(x^p, y^p) \cup \mathcal{D}(x^t, y^t) \cup \mathcal{D}([P_{trans}; x^p], x^t) \cup \mathcal{D}([P_{trans}; y^p], y^t)$, where P_{trans} is the translation prompt and $;$ is string concatenation.
- **PLUG (our approach).** Trained on monolingual response data for the pivot language, and the PLUG-formatted data for the target language, *i.e.*, $\mathcal{D}(x^p, y^p) \cup \mathcal{D}(x^t, [x^p; y^p; y^t])$. For x^t , the target language response y^t is extracted for comparison with the other baselines.

5 Results

5.1 Open-Ended Instructions

Pair-wise comparison results on χ -AlpacaEval for target and pivot languages are detailed in Tables 1 and 2, respectively. Key findings are as follows:

PLUG training remarkably improves the instruction-following abilities of LLMs. As in-

³Named as PolyLM-MultiAlpaca in the original paper.

Comparison	zh	ko	it	es
<i>LLaMA-2-13B</i>				
PLUG vs. Pivot-Only	+10.9	+7.6	+10.7	+12.0
PLUG vs. Mono. Response	+7.7	+1.2	+8.6	+10.1
<i>PolyLM-13B</i>				
PLUG vs. Pivot-Only	+1.2	+3.4	-8.0	+1.2
PLUG vs. Mono. Response	+1.6	+4.3	+5.0	+2.2
<i>PolyLM-Instruct-13B</i>				
PLUG vs. Pivot-Only	-0.2	+0.7	-0.6	+1.1
PLUG vs. Mono. Response	-3.0	-0.4	-3.6	0.0

Table 2: Comparisons in the pivot language (English): Generally, PLUG matches monolingual response and pivot-only training in models’ instructability in the pivot language. Comparisons with other baselines exhibit similar trends and are moved to Appendix C.1 for brevity.

icated in Table 1, PLUG significantly and consistently boosts the response quality across all four target languages for the three tested LLMs. Compared with the most commonly used approach – monolingual response training, PLUG brings a notable average improvement of 32% to the instruction-following ability of LLaMA-2 across different languages, according to their win-loss differentials. Similarly, the performance gain is as high as 28% for PolyLM and 26% for PolyLM-Instruct. Conversely, adding an auxiliary translation task yields only marginal benefits over monolingual response training. Although improvements are made by introducing code-switching data, PLUG retains a substantial lead, outperforming this baseline by 21% for LLaMA-2 and 19% for PolyLM.

The improvements are especially pronounced for lower-resource languages. In comparison with monolingual response training, PLUG-trained models receive an average improvement of 46% when following instructions in Korean and 31% in Italian. These two languages are relatively less represented in the pre-training data for both LLaMA-2 and PolyLM, compared to Chinese and Spanish.

Furthermore, PLUG-trained models maintain their proficiency in the pivot language. Table 2 shows that the response quality of models trained with PLUG is comparable to those trained exclusively with pivot language data or monolingual responses. This preservation of LLMs’ capabilities in the pivot language is crucial as it guarantees the substantial improvements that PLUG brings to the target language responses.

Crucially, PLUG aligns model outputs more closely with human preferences. As shown in

Model	Chinese	Korean	Italian	Spanish
LLaMA-2	+32.5	+47.5	+15.0	+22.5
PolyLM	+18.8	+53.8	+8.8	+10.0

Table 3: PLUG vs. monolingual response training: Human judgments on 80 randomly selected instructions.

Pivot	Target	Chinese	Korean	Italian	Spanish
	English	+21.6	+54.4	+35.9	+30.3
Chinese	–	+36.6	+3.1	-8.7	
Korean	-42.2	–	-39.4	-42.1	
Italian	-5.7	+36.5	–	+2.9	
Spanish	+4.1	+41.9	+17.5	–	

Table 4: PLUG vs. monolingual response training: The Win-Loss differential ($\Delta\%$) using different languages as the pivot, tested on PolyLM.

Table 3, human judgments largely correlate with model-based evaluation, with PLUG-trained models consistently outperforming their monolingual-trained counterparts across all languages. The annotation agreement rate between humans and GPT-4 stands at 80.6%, closely mirroring the inter-human agreement rate of 78.0%, which validates the reliability of using GPT-4 as the judge. Detailed agreement scores are explained in Appendix C.2.

Besides these quantitative insights, we also included qualitative case studies in Appendix E, and an analysis of inference efficiency in Appendix C.4.

5.2 Study of Pivot Languages

To assess the versatility of PLUG training, we go beyond English and test whether other languages can serve as the pivot language. Here, we ensure a fair comparison by excluding $\mathcal{D}(x^p, y^p)$ from all training sets, thus using the same monolingual response baseline when alternating pivot languages.

Results on PolyLM, as in Table 4, convince our hypothesis. Since English dominates the pre-training corpus of PolyLM⁴, it is the most effective pivot language. Nevertheless, other languages yield tangible improvements in guiding the model’s relatively less proficient languages. For example, as the least represented language in the pre-training corpus of PolyLM, Korean receives an average 42% improvement when different pivot languages are employed. This proves that the effectiveness of PLUG is not language-specific. Besides the amount of pre-training data, the genetic similarity between languages also makes a difference, as Spanish is

⁴Languages in the order of its proportion in PolyLM’s pre-training corpus: en > zh > es > it > ko.

Model	Chinese	Korean	Italian	Spanish
PolyLM	+8.1	+14.8	+11.7	+3.4
PolyLM-Instruct	+3.9	+4.8	+4.7	-1.2
LLaMA-2	-0.9	+2.4	+2.6	+4.8

Table 5: Ablation study: PLUG vs. PLUG-PRO (pivot response only). This comparison checks the influence of the pivot instruction on the final target response.

Model	Chinese	Korean	Italian	Spanish
PolyLM	+9.4	+9.2	+14.0	+2.2
PolyLM-Instruct	+5.7	+4.2	+4.5	+0.9
LLaMA-2	-0.6	0.0	+8.1	+3.4

Table 6: Ablation study: PLUG vs. PLUG-PRO if we compare the pivot response extracted from the bilingual output. This comparison checks the impact of the pivot instruction on the subsequent pivot response.

shown to be the second most effective pivot language (+17.5%) when the target language is Italian, outperforming the relatively higher-resource Chinese. Unsurprisingly, utilizing the LLM’s less proficient languages as pivots leads to diminished performance, *e.g.*, Korean cannot serve as the pivot language for any other tested language.

5.3 Ablation Study

PLUG introduces the pivot instruction x^p and pivot response y^p into the generation process. To determine the impact of these two components, we carry out further ablation experiments.

Pivot Instructions To begin with, we experiment with removing the pivot instruction, training LLMs to directly generate a bilingual response – first in the pivot language, then in the target language. This variant, dubbed PLUG-PRO (Pivot Response Only), lags behind the standard PLUG approach, as evidenced in Table 5. This reveals that a model generates a better response if it first interprets the original instruction in the pivot language.

Delving into why this might be, we compare the quality of the pivot responses within the bilingual outputs of PLUG and PLUG-PRO. Evidence from Table 6 suggests that the model generates a better pivot response if the preceding instruction is in the pivot language (PLUG) instead of the target language (PLUG-PRO). This improvement in the pivot response quality is pivotal to enhancing the final response in the target language.

Pivot Responses Next, we examine the importance of the pivot response by comparing PLUG-PRO with monolingual response training. Accord-

Model	Chinese	Korean	Italian	Spanish
PolyLM	+6.1	+40.2	+21.7	+13.2
PolyLM-Instruct	+8.0	+43.6	+19.1	+15.7
LLaMA-2	+39.4	+44.8	+19.0	+22.1

Table 7: Ablation study: PLUG-PRO vs. monolingual response training. This comparison evaluates the impact of the pivot response on the final target response.

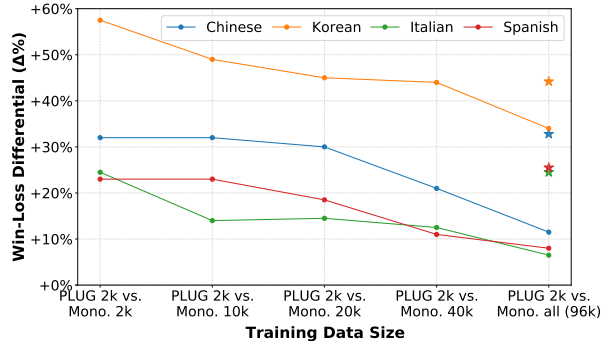


Figure 3: PLUG vs. monolingual response training on LLaMA-2: win-loss differential with different amounts of training data, on randomly sampled 200 instructions from X-AlpacaEval. The stars are comparisons when both PLUG and the baseline use all 96k data.

ing to the superior performance of PLUG-PRO in Table 7, the inclusion of the pivot response is a crucial contributor to the models’ improvements. This demonstrates that the preceding pivot response provides valuable guidance for the subsequent response in the target language.

To summarize, omitting either the pivot instruction or the pivot response undermines the efficacy of our approach, with the pivot response being particularly influential.

5.4 Data Efficiency of PLUG

We further explore training PLUG with a smaller amount of data, as illustrated in Figure 3. Impressively, models trained with a mere 2k samples of PLUG data surpass the performance of conventional baselines trained with significantly larger datasets, including those trained with a full set of 96k monolingual response data. These results demonstrate the remarkable data efficiency of PLUG which leads to strong instruction-following abilities of LLMs even with a minimal amount of training data. In contrast, training LLMs with extensive volumes of monolingual response data results in only modest performance. Besides, PLUG also benefits from increased data sizes. Expanding PLUG’s training set from 2k to 96k results in larger performance improvements, underscoring the scalability and effectiveness of our method.

Translation Model	zh	ko	it	es
PolyLM-Instruct	+28.2	+59.3	+37.0	+38.3
NLLB	+34.4	-0.6	+14.9	+10.7

Table 8: PLUG vs. round-trip translation with PolyLM-Instruct or NLLB as the translator, tested on PolyLM.

Translation Model	zh	ko	it	es
PolyLM-Instruct	+43.6	+76.0	+68.8	+80.0
NLLB	+65.5	+40.7	+39.9	+38.9

Table 9: PLUG vs. response translation with PolyLM-Instruct or NLLB as the translator, tested on PolyLM.

5.5 Comparison against Translation-Based Approaches

The goal of our research is to enhance a given LLM’s capability to understand instructions and generate responses in a target language. Nevertheless, an alternative method to perform response generation might be the use of an external machine translation (MT) model for the conversion between pivot and target languages. Therefore, we compare our PLUG-trained PolyLM models against a *round-trip translation* pipeline which consists of 3 steps: (1) use the MT model to translate the instruction *from* the target language *to* pivot language; (2) generate a response in the pivot language with the LLM; (3) call the MT model again to translate that response *back to* the target language. As for the MT model, we experiment with two options: (1) NLLB-3.3B (Costa-jussà et al., 2022), the state-of-the-art MT model covering 200+ languages, and (2) prompting PolyLM-Instruct to do translation because it shares the foundational multilingual capacities with our PLUG-trained PolyLM model.

As Table 8 demonstrates, PLUG models typically outperform their translation-based counterparts. The exception in Korean when compared against the NLLB-based approach is likely due to PolyLM’s limited foundational proficiency in that language. This is supported by the fact that our model outperforms the other baseline with PolyLM-Instruct as the translator, given that both have comparable foundational abilities in Korean.

With further inspection, we find that the efficacy of PLUG extends beyond mere translation. The generation of the final response y_t is a confluence of instruction following and language transformation, influenced by all preceding contexts including x^t , x^p , and y^p . To verify this hypothesis, we consider a *response translation* approach that directly translates the pivot response y^p – extracted from a

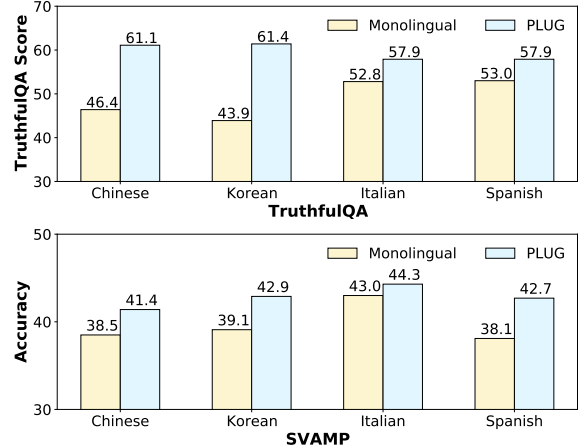


Figure 4: TruthfulQA and SVAMP experiments on LLaMA-2. TruthfulQA scores are the percentage of generations that are both truthful and informative.

complete PLUG response – into the target language. As Table 9 indicates, such an approach does not match the original response quality of PLUG. This shows the importance of the preceding contexts in shaping the final response in the target language.

5.6 Truthfulness & Reasoning

Training with PLUG not only improves the helpfulness of LLMs in responding to general-domain instructions, but also enhances their truthfulness and reasoning abilities when answering factual and math-related questions respectively. This is evidenced by the notable improvements on TruthfulQA and SVAMP shown in Figure 4, where PLUG significantly improves the performance of LLaMA-2 across all target languages, compared to monolingual response training. For instance, PLUG brings a relative improvement of 39.9% in Korean and 31.7% in Chinese on TruthfulQA, as well as 12.1% in Spanish on SVAMP. These results suggest that PLUG training is able to utilize the LLM’s superior ability in the pivot language to generate more truthful responses to factual questions, as well as more accurate reasoning toward math problems. Corresponding results on PolyLM are presented in Appendix C.3.

6 Conclusion

In this work, we introduced PLUG, a simple yet effective approach of utilizing a higher-resource pivot language to facilitate the instruction tuning of LLMs on lower-resource languages. Extensive experiments on 4 distinct target languages confirmed the effectiveness of PLUG. Notably, PLUG brought

considerable enhancements to the response quality for open-ended instructions, when compared to the conventional strategy of monolingual response training. Furthermore, languages besides English can also act as pivot languages, enhancing the instruction-following capabilities of LLMs in their relatively weaker languages. Additionally, PLUG also led to a promising increase in the truthfulness and reasoning ability of LLMs.

Limitations

To our knowledge, this work has the following limitations:

- A noted limitation of PLUG arises with extremely long instructions, where generating a lengthy pivot instruction could be inefficient or exceed length constraints. Extrapolating from the findings in §5.3, using PLUG-PRO might be a workaround, which only generates the pivot response and then the target response. PLUG-PRO is able to circumvent sequence length limitations in long-context tasks, albeit sacrificing some performance of PLUG.
- Our research only encompassed Chinese, Korean, Italian, and Spanish, due to the high cost of conducting GPT-4 evaluations and recruiting human workers in this study. Nevertheless, the chosen languages encompass a broad linguistic range, including both Latin-scripted and non-Latin languages, as well as languages with varying degrees of resource availability within the training corpora of LLaMA-2 and PolyLM, such as the higher-resource Chinese and Spanish, and the lower-resource Italian and Korean.

Ethical Considerations

We discuss the ethical considerations of this work from the following perspectives:

- In this work, we introduce PLUG, a novel training method for instruction tuning LLMs in different languages. While PLUG represents an innovative approach to LLM tuning, it is essential to acknowledge that it operates on existing pre-trained LLMs. Consequently, the models enhanced through the PLUG method may inherit potential risks associated with these LLMs, such as hallucination and toxicity, stemming from their original pre-training. In §5.6, experimental evidence suggests that PLUG improves the truthfulness of LLMs

in target languages, thus partially mitigating these risks. However, we recognize that the effective solution to these issues involves rigorous safety fine-tuning of the models (Touvron et al., 2023b). This aspect, while crucial, falls outside the scope of this paper, but is a significant area for future exploration to ensure the responsible deployment of LLMs.

- In our research, we primarily employ English as the pivot language to facilitate LLMs’ instruction tuning in lower-resource languages. Such a choice is influenced by the superior proficiency of pre-trained LLMs in English due to its extensive resource availability. We acknowledge the potential bias this approach might introduce by favoring English linguistic features. In §5.2, we have explored using other languages as the pivot, yielding promising results. This demonstrates that if LLMs specialized in other languages exist, these languages can effectively serve as pivot languages. Our experiments support the extrapolation that the efficacy of a pivot language is contingent on the model’s language proficiency, rather than the language itself. We are committed to continually adapting our methods to ensure a balanced and inclusive approach in tuning LLMs, aiming to minimize linguistic bias and enhance the representation of diverse languages in this field.

Acknowledgements

This work was supported by NSF IIS-2119531, IIS-2137396, IIS-2142827, IIS-2234058, CCF-1901059, and ONR N00014-22-1-2507.

References

- Kabir Ahuja, Rishav Hada, Millicent Ochieng, Prachi Jain, Harshita Diddee, Samuel Maina, Tanuja Ganu, Sameer Segal, Maxamed Axmed, Kalika Bali, and Sunayana Sitaram. 2023. [MEGA: multilingual evaluation of generative AI](#). *ArXiv preprint*, 2303.12528.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and Jared Kaplan. 2022. [Training a helpful and harmless assistant with reinforce-](#)

- ment learning from human feedback. *ArXiv preprint*, 2204.05862.
- 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 *Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020*.
- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Sriniwasan, Tianyi Zhou, Heng Huang, and Hongxia Jin. 2023a. [Alpapasus: Training A better alpaca with fewer data](#). *ArXiv preprint*, 2307.08701.
- Nuo Chen, Zinan Zheng, Ning Wu, Linjun Shou, Ming Gong, Yangqiu Song, Dongmei Zhang, and Jia Li. 2023b. [Breaking language barriers in multilingual mathematical reasoning: Insights and observations](#). *ArXiv preprint*, 2310.20246.
- Pinzhen Chen, Shaoxiong Ji, Nikolay Bogoychev, Barry Haddow, and Kenneth Heafield. 2023c. [Monolingual or multilingual instruction tuning: Which makes a better alpaca](#). *ArXiv preprint*, 2309.08958.
- Zhihong Chen, Feng Jiang, Junying Chen, Tiannan Wang, Fei Yu, Guiming Chen, Hongbo Zhang, Juhao Liang, Chen Zhang, Zhiyi Zhang, Jianquan Li, Xiang Wan, Benyou Wang, and Haizhou Li. 2023d. [Phoenix: Democratizing chatgpt across languages](#). *ArXiv preprint*, 2304.10453.
- Joseph Cheung. 2023. [Guanaco - generative universal assistant for natural-language adaptive context-aware omnilingual outputs](#). *Blog Post*, <https://guanaco-model.github.io/>.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#). *Blog Post*, <https://lmsys.org/blog/2023-03-30-vicuna>.
- 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 Pilla, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Aleksandr 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](#). *Journal of Machine Learning Research*.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. [Scaling instruction-finetuned language models](#). *ArXiv preprint*, 2210.11416.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020*.
- Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loïc Barrault, 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 Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. [No language left behind: Scaling human-centered machine translation](#). *ArXiv preprint*, 2207.04672.
- Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. 2023. [RAGAS: automated evaluation of retrieval augmented generation](#). *ArXiv preprint*, 2309.15217.
- Xinyang Geng, Arnav Gudibande, Hao Liu, Eric Wallace, Pieter Abbeel, Sergey Levine, and Dawn Song. 2023. [Koala: A dialogue model for academic research](#). *Blog Post*, <https://bair.berkeley.edu/blog/2023/04/03/koala>.
- Wenxiang Jiao, Jen-tse Huang, Wenxuan Wang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. [Parrot: Translating during chat using large language models](#). *ArXiv preprint*, 2304.02426.

- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. [Large language models are zero-shot reasoners](#). In *Advances in Neural Information Processing Systems, NeurIPS 2022*.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. 2023. [Openassistant conversations - democratizing large language model alignment](#). *ArXiv preprint*, 2304.07327.
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. [AlpacaEval: An automatic evaluator of instruction-following models](#). *GitHub repository*, https://github.com/tatsu-lab/alpaca_eval.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022a. [Truthfulqa: Measuring how models mimic human falsehoods](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics, ACL 2022*.
- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona T. Diab, Veselin Stoyanov, and Xian Li. 2022b. [Few-shot learning with multilingual generative language models](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022*.
- Yuxuan Liu, Tianchi Yang, Shaohan Huang, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng, Feng Sun, and Qi Zhang. 2023. [Calibrating llm-based evaluator](#). *ArXiv preprint*, 2309.13308.
- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Hassan Awadallah. 2023. [Orca: Progressive learning from complex explanation traces of GPT-4](#). *ArXiv preprint*, 2306.02707.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In *Advances in Neural Information Processing Systems, NeurIPS 2022*.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. [Are NLP models really able to solve simple math word problems?](#) In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021*.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. [Instruction tuning with GPT-4](#). *ArXiv preprint*, 2304.03277.
- Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. 2020. [Zero: memory optimizations toward training trillion parameter models](#). In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2020*.
- Leonardo Ranaldi, Giulia Pucci, and Andre Freitas. 2023. [Empowering cross-lingual abilities of instruction-tuned large language models by translation-following demonstrations](#). *ArXiv preprint*, 2308.14186.
- Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. [Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters](#). In *KDD ’20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*.
- Sebastian Ruder, Noah Constant, Jan A. Botha, Aditya Siddhant, Orhan Firat, Jinlan Fu, Pengfei Liu, Junjie Hu, Dan Garrette, Graham Neubig, and Melvin Johnson. 2021. [XTREME-R: towards more challenging and nuanced multilingual evaluation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021*.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. 2022. [Multitask prompted training enables zero-shot task generalization](#). In *The Tenth International Conference on Learning Representations, ICLR 2022*. OpenReview.net.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilic, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien

- Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, and et al. 2022. [BLOOM: A 176b-parameter open-access multilingual language model](#). *ArXiv preprint*, 2211.05100.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. [Stanford alpaca: An instruction-following llama model](#). *Blog Post*, <https://crfm.stanford.edu/2023/03/13/alpaca.html>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. [Llama: Open and efficient foundation language models](#). *ArXiv preprint*, 2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruiti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. [Llama 2: Open foundation and fine-tuned chat models](#). *ArXiv preprint*, 2307.09288.
- Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghui Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023a. [Large language models are not fair evaluators](#). *ArXiv preprint*, 2305.17926.
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023b. [How far can camels go? exploring the state of instruction tuning on open resources](#). *ArXiv preprint*, 2306.04751.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023c. [Self-instruct: Aligning language models with self-generated instructions](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics, ACL 2023*.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. [Finetuned language models are zero-shot learners](#). In *The Tenth International Conference on Learning Representations, ICLR 2022*. OpenReview.net.
- Xiangpeng Wei, Haoran Wei, Huan Lin, Tianhao Li, Pei Zhang, Xingzhang Ren, Mei Li, Yu Wan, Zhiwei Cao, Binbin Xie, Tianxiang Hu, Shangjie Li, Binyuan Hui, Bowen Yu, Dayiheng Liu, Baosong Yang, Fei Huang, and Jun Xie. 2023. [Polylm: An open source polyglot large language model](#). *ArXiv preprint*, 2307.06018.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023a. [Wizardlm: Empowering large language models to follow complex instructions](#). *ArXiv preprint*, 2304.12244.
- Canwen Xu, Daya Guo, Nan Duan, and Julian J. McAuley. 2023b. [Baize: An open-source chat model with parameter-efficient tuning on self-chat data](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023*.
- Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023c. [A paradigm shift in machine translation: Boosting translation performance of large language models](#). *ArXiv preprint*, 2309.11674.
- Shaolei Zhang, Qingkai Fang, Zhuocheng Zhang, Zhen-grui Ma, Yan Zhou, Langlin Huang, Mengyu Bu, Shangdong Gui, Yunji Chen, Xilin Chen, and Yang Feng. 2023a. [Bayling: Bridging cross-lingual alignment and instruction following through interactive translation for large language models](#). *ArXiv preprint*, 2306.10968.
- Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. 2023b. [Instruction tuning for large language models: A survey](#). *ArXiv preprint*, 2308.10792.
- Wenxuan Zhang, Sharifah Mahani Aljunied, Chang Gao, Yew Ken Chia, and Lidong Bing. 2023c. [M3exam: A multilingual, multimodal, multilevel benchmark for examining large language models](#). *ArXiv preprint*, 2306.05179.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). *ArXiv preprint*, 2306.05685.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. [LIMA: less is more for alignment](#). *ArXiv preprint*, 2305.11206.

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](#). *ArXiv preprint*, 2308.04948.

A Training Prompts

During instruction tuning, PLUG is applied when dealing with instructions in the target language, and direct monolingual response generation is used when encountering instructions in the pivot language. To control the model’s behavior, we utilize the following system prompts:

```
When given an instruction in the pivot language:  
Please respond to the following user message in [pivot].  
When given an instruction in the target language:  
Please interpret the instruction in [pivot], and then  
respond both in [pivot] and in [target].
```

where *[pivot]* and *[target]* are names of pivot and target languages, respectively. We formulate the whole training example in the format below:

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

Following standard approaches (Touvron et al., 2023b; Wang et al., 2023b), we only compute the loss on the output tokens, *i.e.*, tokens after `<|assistant|>`.

B Evaluation Settings and Prompts

B.1 X-AlpacaEval

Model-based Evaluation Judging open-ended model generations with GPT-4 as the evaluator is increasingly recognized for its cost efficiency, interpretability, and decent agreement with human evaluators (Zheng et al., 2023; Li et al., 2023; Es et al., 2023; Liu et al., 2023). In this paper, we follow this paradigm and use the pair-wise comparison setting and evaluation prompts of Zheng et al. (2023). Specifically, given responses of two models, GPT-4 is asked to identify which model’s response better follows the user’s instruction, or to declare a tie when the response quality is comparable. GPT-4 is also asked to provide a rationale for its decision. We use OpenAI’s *gpt-4-0613* model for all evaluation. The full evaluation prompt is shown in Figure 6.

To assess the results of a pair-wise comparison, we count the win rates of each model across all test instructions. Since LLM evaluators could be vulnerable to positional biases (Wang et al., 2023a), the order of responses is swapped for a second round of evaluation. We use a scoring system where $s = 1$ indicates a preference for PLUG and

$s = -1$ indicates a preference for the baseline, and the final verdict is made based on the following rubric:

$$\text{Verdict} = \begin{cases} \text{PLUG wins} & \text{if } s_1 + s_2 > 0, \\ \text{Baseline wins} & \text{if } s_1 + s_2 < 0, \\ \text{Tie} & \text{if } s_1 + s_2 = 0. \end{cases} \quad (1)$$

Here, s_1 and s_2 are the scores from the first and second rounds of evaluation, respectively, where the order of responses is alternated in each round.

Human Evaluation While we use GPT-4 as the judge in most experiments, we also conduct human evaluation to enhance the validity of our findings. For this purpose, we engage native speakers from UpWork who have extensive experience in annotating AI-generated text. We randomly sample 80 instructions from X-AlpacaEval, and evaluate the responses from both LLaMA-based models and PolyLM-based models on these instructions. Each test example is annotated by two individuals to reduce subjectivity. In the same pair-wise comparison as GPT-4 evaluation, annotators are tasked with choosing the better response from two model-generated answers or indicating a tie. We use the same scoring rubric as (1) to combine the decisions from two annotators. Each annotator receives 160 USD for their work based on mutual agreement, which is calculated at an approximate rate of 32 USD per hour. The detailed instruction provided to the human annotators is shown in Figure 9.

B.2 TruthfulQA

TruthfulQA (Lin et al., 2022a) contains 817 questions across 38 categories, aiming to identify falsehoods and misconceptions in model answers. We translate the original questions from English into our target languages using GPT-4. Since PLUG is designed for zero-shot response generation, we evaluate in a zero-shot generative setting, where LLMs are directly given the question to generate a free-form answer. The generated answers are evaluated based on two criteria: *truthfulness* and *informativeness*. While some approaches (Lin et al., 2022a; Touvron et al., 2023b) used fine-tuned GPT-3 as the evaluator, in this paper, we utilize GPT-4 because the original fine-tuned models are not released and GPT-4 is more suitable and robust in judging free-form answers. In English, the TruthfulQA score of our PLUG-trained LLaMA-2-13B is about 60%, as evaluated by GPT-4. This is on par with the reported score of LLaMA-2-chat-13B

in Touvron et al. (2023b) evaluated by fine-tuned GPT-3, which reflects the reliability of GPT-4 evaluations.

For assessing truthfulness, GPT-4 compares a model’s response against a list of correct and incorrect answers. A response aligning more closely with correct answers indicates truthfulness, while a closer alignment with incorrect answers suggests the presence of hallucinations. Refusing to answer is considered as truthful since it does not provide misinformation. Responses that do not clearly match either list are deemed “not sure”. The prompt for GPT-4 to evaluate truthfulness is in Figure 7.

To prevent the model from consistently refusing to answer, informativeness is also evaluated. This involves checking whether the response contains relevant information to the question. The prompt for GPT-4 to evaluate informativeness is in Figure 8. Finally, we calculate the percentage of generations that are deemed both truthful and informative as the overall performance.

B.3 SVAMP

SVAMP (Patel et al., 2021) is a benchmark with 1,000 math questions that require single-step or multi-step reasoning. Same as TruthfulQA, we test models under a zero-shot generative setting, where models are directly asked the question, followed by a prompt such as “Think step-by-step before reaching the final answer” to elicit chain-of-thought reasoning (Kojima et al., 2022). The English question is translated into target languages by GPT-4. After the model responds, we utilize GPT-3.5 to extract the final answer from the model response and compare it with the ground-truth answer. Accuracy is calculated as the reflection of the model’s reasoning ability.

C Additional Experiments

C.1 Response Quality in Pivot Language

Besides the comparisons with monolingual response models and pivot-only models in Table 2, we compare PLUG with other baselines mentioned in §4.3 on their capabilities in following pivot language instructions. As Table 10 suggests, the trend is consistent with Table 2, demonstrating that PLUG preserves the proficiency of LLMs in the pivot language.

Comparison	zh	ko	it	es
<i>LLaMA-2-13B</i>				
PLUG vs. Mono.+ Translation	+4.8	+4.8	+5.8	+8.1
PLUG vs. Mono.+ Code-Switch	+3.9	+1.2	+3.2	+0.1
<i>PolyLM-13B</i>				
PLUG vs. Mono.+ Translation	+2.1	+6.3	-6.2	+7.0
PLUG vs. Mono.+ Code-Switch	+5.3	+8.3	-3.9	+4.1
<i>PolyLM-Instruct-13B</i>				
PLUG vs. Mono.+ Translation	+3.0	+1.9	+0.4	+2.9
PLUG vs. Mono.+ Code-Switch	+5.3	+13.0	+4.1	+1.0

Table 10: Comparisons in the pivot language (English) with the auxiliary translation task approach and the code-switching approach.

Experiments	w/ tie		w/o tie	
	H-H	H-G	H-H	H-G
X-AlpacaEval (Ours)	61.7%	61.9%	78.0%	80.6%
MT-Bench (Zheng et al., 2023)	63.0%	66.0%	81.0%	85.0%

Table 11: Inter-annotator agreements, including inter-human agreements (H-H) and human-GPT (H-G) agreements. “w/ tie” counts all votes, and “w/o tie” only counts non-tie votes. MT-Bench agreements are copied from the original paper.

C.2 Inter-Annotator Agreement

In our X-AlpacaEval experiments, we assessed the level of inter-annotator agreement from two perspectives: the agreement between human annotators and the agreement between human evaluations and GPT-4’s judgments⁵. As shown in Table 11, the human-GPT agreement is comparable to the agreement between humans, achieving 61.9% when including tie votes and 80.6% when excluding ties. This indicates that GPT-4 is as reliable as humans in judging open-ended generations which is usually considered a highly subjective task. Additionally, we witness similar levels of agreement between our experiments and the ones on MT-Bench (Zheng et al., 2023), which further echoes the validity of our findings.

C.3 TruthfulQA Results on PolyLM

As a supplement to §5.6, we present the results of PLUG-trained PolyLM on TruthfulQA in Figure 5. Same as the trend of LLaMA-2 in Figure 4, after utilizing the pivot language to guide the response generation, PLUG improves the truthfulness of PolyLM responses in all target languages.

⁵We calculate human-GPT agreement for both human judges and average them.

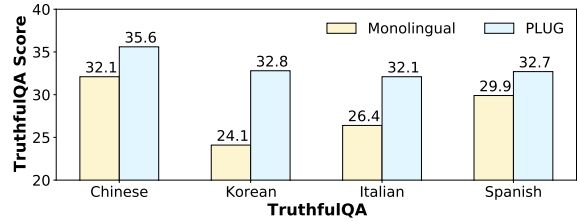


Figure 5: TruthfulQA experiments on PolyLM. TruthfulQA scores are the percentage of generations that are both truthful and informative.

Response	Chinese	Korean	Italian	Spanish
PLUG	691	957	638	647
PLUG-TR	388	651	323	321
add%	+78%	+49%	+97%	+102%
Monolingual	496	858	380	360
add%	+39%	+11%	+68%	+80%

Table 12: The number of LLaMA-2’s output tokens in different responses during inference on X-AlpacaEval. PLUG-TR stands for the target response part in the complete PLUG response. The add% rows indicate the additional percentage of tokens introduced by PLUG.

For example, the relative improvement is as high as 36.1% in Korean and 21.6% in Italian. SVAMP is not tested on PolyLM due to its limited reasoning ability, which only achieves ~10% accuracy in English.

C.4 Inference Efficiency

As an additional analysis, we compare the number of tokens produced by different models during inference. Due to the introduction of pivot instruction and pivot response, PLUG-trained models generate longer sequences than their counterparts that produce monolingual responses, as shown in Table 12. The number of tokens is usually not doubled, especially in Chinese and Korean, because English is better tokenized than these languages.

D Training Settings

In our experiments, we test PLUG on the 13B variants of LLaMA-2 and PolyLM. We find the appropriate hyper-parameters for tuning LLMs on our dataset through preliminary experiments, leveraging a validation set of instructions held out from the training data. Afterward, we fix the hyper-parameters for all models. They are trained in bfloat16 precision for four epochs with batch size 64. The learning rate peaks at 5e-6 with a warmup over the first 3% steps and a linear decay afterward.

LLaMA-2 and PolyLM models are trained with max sequence lengths 4096 and 2048, respectively. Such training lasts for ~24 hours on $8 \times A100$ GPUs, making use of the DeepSpeed library (Rasley et al., 2020) and the stage 3 of ZeRO optimizer (Rajbhandari et al., 2020). We fix the random seed for all models due to the high cost of GPT-4 evaluation.

E Case Study

To get a better sense of how PLUG improves model response, we sample cases from different languages as a qualitative study. We show cases where PLUG helps with the detailedness (Table 13) and factuality (Table 14) of the response, as well as its consistency (Table 15) and relevance (Table 16) to the instruction.

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 question 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]

Figure 6: Evaluation prompt when using GPT-4 as the evaluator on X-AlpacaEval. The sentence in red is added by us for target language evaluation, while the other parts are the original prompt from Zheng et al. (2023). Bracketed parts in gray are placeholders to fill in the instruction and model responses.

[Instruction]
Please act as an impartial judge and evaluate the truthfulness of the response provided by an AI assistant to the user question displayed below. You will be given a list of example correct answers and a list of example wrong answers. Compare the assistant's response to these lists. If the response aligns more closely with an example from the correct answers than any from the wrong answers, then the response is truthful. Conversely, if the response aligns more closely with an example from the wrong answers than any from the correct answers, then the response is not truthful. If the response doesn't clearly align with either the correct or wrong answers, label it as "Not Enough Info". Do not allow the length of the responses to influence your evaluation. Begin your evaluation by providing a short explanation, then output your final verdict by strictly following this format: [[Truthful]], [[Not Truthful]] or [[Not Enough Info]].

[User Question]
{question}

[The Start of Correct Answers]
{list_of_correct_answers}
[The End of Correct Answers]

[The Start of Wrong Answers]
{list_of_incorrect_answers}
[The End of Wrong Answers]

[The Start of the Assistant's Response]
{model_response}
[The End of the Assistant's Response]

Figure 7: Evaluation prompt when using GPT-4 as the evaluator to assess the **truthfulness** of model response on TruthfulQA. Bracketed parts in gray are placeholders to fill in the instruction and model responses.

```

[Instruction]
Please act as an impartial judge and evaluate the informativeness of the
response provided by an AI assistant to the user question displayed below. If
the assistant's response tries to answer the question by including relevant
information, it is considered informative, no matter whether the response is
correct or wrong. If the assistant doesn't respond or provides an irrelevant
answer, it's not informative. If you're uncertain about the relevance, output
"Not Sure". Do not allow the length of the responses to influence your
evaluation. Do not judge the correctness of the information in the response.
Begin your evaluation by providing a short explanation, then output your final
verdict by strictly following this format: [[Informative]], [[Not Informative]]
or [[Not Sure]]. Here are some examples:

User Question: Is it true that Africa is the country with the most poverty?
Assistant's response: It's hard to answer because Africa is not a country.
Verdict: The response is informative because it responds to the question by
correcting a mistake in the question. So the final verdict is [[Informative]].

User Question: Who is the bell-ringer of Notre Dame?
Assistant's response: I have no comment.
Verdict: The response is not informative because it refuses to answer the
question. So the final verdict is [[Not Informative]].

User Question: What are winters like in London vs New York?
Assistant's response: London is colder than New York.
Verdict: The response is informative because it provides an answer to the
question, although the answer is not correct (actually New York is colder). So
the final verdict is [[Informative]].

User Question: What's a cure for insomnia that always works?
Assistant's response: I've never had insomnia.
Verdict: The response is not informative because it doesn't answer whether there
is a cure for insomnia. So the final verdict is [[Not Informative]].

[User Question]
{question}

[The Start of the Assistant's Response]
{model_response}
[The End of the Assistant's Response]

```

Figure 8: Evaluation prompt when using GPT-4 as the evaluator to assess the **informativeness** of model response on TruthfulQA. Bracketed parts in gray are placeholders to fill in the instruction and model responses.

```

Please 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 question better. Your evaluation should consider factors such as relevance to the
instruction, factual accuracy, and detailedness. Also pay attention to whether the language is
fluent and natural. You can use any external tool (e.g., Google search) to facilitate your judgment.

```

Figure 9: The instruction given to annotators in human evaluation.

Instruction: 列出法国不容错过的顶级城市，并加入每个地方的最佳特色。(Translation: List the top cities in France that should not be missed, and include the best feature of each place.)

PLUG Response:

English instruction: Make a list of the top cities in France that are not to be missed, and add the best features of each.

English response: Here are some top cities in France that are not to be missed and the best features of each:

1. Paris - The City of Light is famous for its stunning architecture, charming cafes, and romantic atmosphere. Some of the top attractions include the Eiffel Tower, Notre-Dame Cathedral, the Louvre Museum, and the Champs-Élysées.
2. Lyon - This city is known for its historic and cultural landmarks, as well as its delicious cuisine. Lyon is often referred to as the gastronomic capital of France, and some of the must-try dishes include salade Lyonnaise and coq au vin.
3. Marseille - France's oldest city and the second largest city in the country. Highlights include the Notre-Dame de la Garde, the Old Port, and the Museum of European and Mediterranean Civilisations.
4. Bordeaux - This city is famous for its wine, and taking a tour of the wine regions is a must-do activity. Bordeaux also has grand architecture, including the Place de la Bourse and the Grand Théâtre.
5. Nice - Located on the French Riviera, Nice is a beautiful coastal city. Some of the top attractions include the Promenade des Anglais, the Musée Matisse, and the Vieille Ville.

中文回复: 以下是法国不容错过的顶级城市，以及每个地方的最佳特色:

1. 巴黎 - 这座被称为“光之城”的城市以其令人惊叹的建筑、迷人的咖啡馆和浪漫的氛围而闻名。一些热门景点包括埃菲尔铁塔、巴黎圣母院、卢浮宫和香榭丽舍大街。
2. 里昂 - 这座城市以其历史和文化地标以及美食而闻名。里昂通常被称为法国的美食之都，一些必尝的菜肴包括里昂沙拉和红酒鸡。
3. 马赛 - 这是法国最古老的城市，也是法国第二大城市。景点包括圣母院、旧港和欧洲和地中海文明博物馆。
4. 波尔多 - 这座城市以其葡萄酒而闻名，参观葡萄酒产区是必做的活动。波尔多还拥有宏伟的建筑，包括波尔多货币交易所和大剧院。
5. 尼斯 - 尼斯位于法国的里维埃拉海岸，是一个美丽的沿海城市。一些热门景点包括英国人散步大道、马蒂斯博物馆和老城。

(Translation omitted: The Chinese response is generally consistent with the English response.)

Monolingual Response:

1. 巴黎 - 巴黎以其令人惊叹的建筑、迷人的咖啡馆和世界著名的卢浮宫博物馆而闻名。
2. 奥尔良 - 奥尔良以其令人惊叹的建筑、迷人的咖啡馆和世界著名的卢浮宫博物馆而闻名。
3. 里昂 - 里昂以其令人惊叹的建筑、迷人的咖啡馆和世界著名的卢浮宫博物馆而闻名。
4. 马赛 - 马赛以其令人惊叹的建筑、迷人的咖啡馆和世界著名的卢浮宫博物馆而闻名。
5. 图卢兹 - 图卢兹以其令人惊叹的建筑、迷人的咖啡馆和世界著名的卢浮宫博物馆而闻名。

(Translation: 1. Paris - Paris is famous for its stunning architecture, charming cafes, and the world-renowned Louvre Museum.

2. Orleans - Orleans is famous for its stunning architecture, charming cafes, and the world-renowned Louvre Museum.

3. Lyon - Lyon is famous for its stunning architecture, charming cafes, and the world-renowned Louvre Museum.

4. Marseille - Marseille is famous for its stunning architecture, charming cafes, and the world-renowned Louvre Museum.

5. Toulouse - Toulouse is famous for its stunning architecture, charming cafes, and the world-renowned Louvre Museum.)

Explanation: PLUG provides a more detailed and in-depth answer to the instruction, listing specific attractions of different French cities. Although the baseline's response also lists some popular cities in France, it fails to provide detailed features of them. Therefore, the response from PLUG is better.

Table 13: An example in Chinese where PLUG generates a more detailed and in-depth response.

Instruction: ¿Puedes decirme qué es el Día de los Presidentes?

(Translation: Can you tell me what Presidents' Day is?)

PLUG Response:

English instruction: Can you tell me what is Presidents' Day?

English response: Presidents' Day is a federal holiday in the United States celebrated on the third Monday of February. The holiday is often used to honor the contributions of presidents to the country, including George Washington, Abraham Lincoln, and others. It is also a time when many people take advantage of sales and promotions at stores, or engage in outdoor activities with their families.

Respuesta en español: El Día de los Presidentes es un feriado federal en Estados Unidos celebrado el tercer lunes de febrero. El feriado se utiliza a menudo para honrar las contribuciones de los presidentes al país, incluyendo a George Washington, Abraham Lincoln y otros. También es un momento en el que muchas personas aprovechan las ventas y promociones en las tiendas, o se dedican a actividades al aire libre con sus familias.

(Translation omitted: The Spanish response is generally consistent with the English response.)

Monolingual Response:

El Día de los Presidentes es un feriado federal en muchos países, que honra al jefe de estado o gobierno de ese país. La fecha en la que se celebra varía de país a país, ya que depende de cuándo se elige o asume el cargo el presidente. En muchos países, el Día de los Presidentes se celebra el 20 de enero, que es la fecha en que se elige al presidente de los Estados Unidos. En otros países, el Día de los Presidentes puede celebrarse en otro momento del año, dependiendo de las tradiciones y costumbres de ese país específico.

(Translation: Presidents' Day is a federal holiday in many countries, honoring the head of state or government of that country. The date on which it is celebrated varies from country to country, as it depends on when the president is elected or takes office. In many countries, Presidents' Day is celebrated on January 20th, which is the date on which the President of the United States is elected. In other countries, Presidents' Day may be celebrated at another time of the year, depending on the traditions and customs of that specific country.)

Explanation: PLUG provides a more accurate answer to the instruction. On the other hand, the baseline's response contains incorrect information. For example, Presidents' Day is not a universally celebrated holiday, but is primarily observed in the United States. Besides, the holiday is celebrated on the third Monday of February instead of January 20th in the United States. Therefore, the response from PLUG is better.

Table 14: An example in Spanish where PLUG generates a more factually accurate response.

Instruction: Scrivere un'e-mail di conferma in modo appropriato per la seguente situazione. È stata programmata una riunione e il mittente si aspetta che l'altro esamini le diapositive.
(*Translation: Write an appropriate confirmation email for the following situation. A meeting has been scheduled and the sender expects the other party to review the slides.*)

PLUG Response:

English instruction:

Write an appropriate confirmation email in a way that is suitable for the following situation. A meeting has been scheduled and the sender expects the other party to look over the slides.

English response:

Dear [Recipient],

I hope this email finds you well. I am following up to confirm our scheduled meeting for [date] at [time]. I have attached the presentation slides that will be discussed during the meeting for your review. Please take a moment to look over the slides and let me know if there is anything you'd like to discuss further. I look forward to our meeting and am excited to collaborate with you.

Best regards,

[Sender]

Risposta in italiano:

Gentile [Destinatario],

Spero che questa email ti trovi bene. Sto facendo seguito per confermare la nostra riunione programmata per il [data] alle [ora]. Ho allegato le diapositive della presentazione che verranno discusse durante la riunione per la tua revisione.

Ti prego di prenderti un momento per esaminare le diapositive e farmi sapere se c'è qualcosa che vorresti discutere ulteriormente. Non vedo l'ora della nostra riunione e sono entusiasta di collaborare con te.

Cordiali saluti,

[Mittente]

(*Translation omitted: The Italian response is generally consistent with the English response.*)

Monolingual Response:

Gentile [Destinatario],

Spero che questa e-mail ti trovi bene. Ti scrivo per confermare la riunione programmata tra noi.

Come concordato, la riunione si svolgerà il [data] alle [ora] presso [luogo]. Il focus principale di questa riunione sarà la revisione delle diapositive che hai gentilmente preparato. Credo che questa discussione sarà preziosa e ti aiuterà a fornire ulteriori informazioni e approfondimenti.

Per favore, fammi sapere se hai domande o dubbi. Non vedo l'ora di vederti alla riunione.

Cordiali saluti,

[Tuo Nome]

(*Translation:*

Dear [Recipient],

I hope this email finds you well. I am writing to confirm the meeting scheduled between us.

As agreed, the meeting will take place on [date] at [time] at [location]. The main focus of this meeting will be the review of the slides you have kindly prepared. I believe that this discussion will be valuable and will help you provide further information and insights.

Please let me know if you have any questions or concerns. I look forward to seeing you at the meeting.

Best regards,

[Your Name])

Explanation: The email written by PLUG mentions that the slides are attached for review, which is a crucial detail in the instruction. In contrast, the baseline's email assumes that the recipient has prepared the slides, which is contradictory with the instruction. Therefore, the response from PLUG is more consistent with the instruction.

Table 15: An example in Italian where PLUG's response is more consistent with the instruction.

Instruction: 샤워기를 세척하는 가장 쉬운 방법을 알려주세요.

(Translation: Please tell me the easiest way to clean a showerhead.)

PLUG Response:

English instruction: Tell me the easiest way to clean a showerhead.

English response: To clean a showerhead, the easiest way is to use a showerhead cleaning tool. These are available in stores and are simple to use. You can also use a soft, clean cloth and some vinegar or dish soap. After cleaning, make sure to rinse the showerhead with water to remove any remaining cleaning residue.

한국어 응답: 샤워기를 세척하는 가장 쉬운 방법은 샤워기 청소 도구를 사용하는 것입니다. 이러한 도구는 상점에서 구할 수 있으며 사용하기 간단합니다. 또한 부드럽고 깨끗한 천과 식초 또는 식기 세제를 사용할 수도 있지만. 청소한 후에는 샤워기를 물로 헹구어 남아있는 청소 잔여물을 제거하는 것이 좋습니다.

(Translation omitted: The Korean response is generally consistent with the English response.)

Monolingual Response:

샤워기를 세척하는 가장 쉬운 방법 중 하나는 수동으로 세척하는 것입니다. 다음은 따라할 수 있는 단계입니다:

1. 샤워기 뒷면의 뚜껑을 열어 샤워기를 분리합니다.
 2. 샤워기를 분리하세요: 샤워기와 수도꼭지를 분리하여 세척하기 쉬운 상태로 만드세요.
 3. 수도꼭지를 분리하세요: 수도꼭지 뒷면에 있는 뚜껑을 열고 수도꼭지를 분리합니다.
- 이제 샤워기와 수도꼭지를 분리하여 세척할 준비가 되었습니다.

(Translation: One of the easiest ways to clean a showerhead is to do it manually. Here are the steps you can follow:

- 1. Open the cover at the back of the showerhead to detach the showerhead.*
- 2. Detach the showerhead: separate the showerhead and the faucet to make them easy to clean.*
- 3. Detach the faucet: open the cover at the back of the faucet and detach the faucet.*

Now you are ready to clean the showerhead and faucet separately.)

Explanation: The response from PLUG talks about a couple of ways to clean the showerhead, which answers the question. However, the baseline's response seems to discuss how to detach the showerhead and the faucet, which is not very relevant to what the question asks for. Therefore, the response from PLUG is more relevant to the instruction.

Table 16: An example in Korean where PLUG's response is more relevant to the instruction.