

LLaMA Beyond English: An Empirical Study on Language Capability Transfer

Jun Zhao*, Zhihao Zhang*, Luhui Gao, Qi Zhang[†], Tao Gui, Xuanjing Huang

¹School of Computer Science, Fudan University
{zhaoj19,zhangzhihao19,qz,tgui}@fudan.edu.cn

Abstract

In recent times, substantial advancements have been witnessed in large language models (LLMs), exemplified by ChatGPT, showcasing remarkable proficiency across a range of complex tasks. However, many mainstream LLMs (e.g. LLaMA) are pretrained on English-dominant corpus, which limits their performance in other non-English languages. In this paper, we focus on how to effectively transfer the capabilities of language generation and following instructions to a non-English language. To answer this question, we conduct an extensive empirical investigation based on LLaMA, accumulating over 1440 GPU hours. We analyze the impact of key factors such as vocabulary extension, further pretraining, and instruction tuning on transfer. To accurately assess the model’s level of knowledge, we employ four widely used standardized testing benchmarks: C-Eval, MMLU, AGI-Eval, and GAOKAO-Bench. Furthermore, a comprehensive evaluation of the model’s response quality is conducted, considering aspects such as accuracy, fluency, informativeness, logical coherence, and harmlessness, based on LLM-Eval, a benchmarks consisting instruction tasks from 17 diverse categories. Our evaluation results demonstrate that comparable performance to state-of-the-art transfer models can be achieved with less than 1% of the pretraining data, both in terms of knowledge alignment and response quality. Furthermore, the experimental outcomes across the thirteen low-resource languages also exhibit similar trends. We anticipate that the conclusions revealed by the experiments will aid the community in developing non-English LLMs.

Introduction

For decades, researchers in Natural Language Processing (NLP) have been exploring the fundamental principles of intelligence (Bubeck et al. 2023). The recent advances in large language models (LLMs) seem to have revealed a glimmer of hope. Benefitting from the unprecedented scales of model size and training data, many LLMs like ChatGPT (OpenAI 2022), PaLM (Anil et al. 2023), LLaMA (Touvron et al. 2023a), and others have emerged strong capabilities in reasoning (Cobbe et al. 2021), planning (Huang et al. 2022), and learning from experience (Dong et al. 2023) at or surpassing human levels. These general capabilities also provide a foundation for LLMs to address intricate

*These authors contributed equally.

[†]Corresponding Author

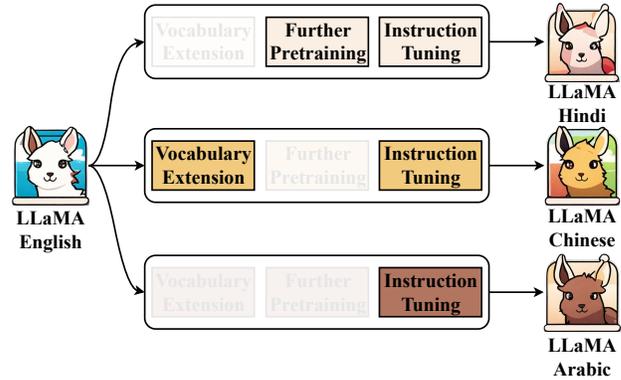


Figure 1: Pretrained LLaMA models, which are primarily trained on English-dominated corpus (as depicted on the left), are not inherently proficient in handling non-English languages. We aim to investigate the necessity of vocabulary extension, further pretraining, and instruction tuning, as well as to what extent they influence the capability transfer. This exploration enables us to efficiently transfer LLaMA’s language capabilities to non-English languages (as illustrated on the right), minimizing costs in the process.

real-world tasks, such as successfully completing the entire Uniform Bar Examination (UBE) (Katz et al. 2023) or coding based on natural language instructions (StabilityAI 2023).

Many well-known LLMs are capable of comprehending input and generating responses across different languages, thanks to their pretraining on a diverse mix of corpus from multiple languages. However, due to the imbalanced distribution of language resources, collecting extensive training data for all languages is nearly impossible (Ranta and Goutte 2021). Taking the representative LLM BLOOM (Scao et al. 2023) as an example, it has been pretrained on 46 natural languages. Yet, this number accounts for only 0.66% of the roughly 7,000 languages currently in use. Moreover, within the corpus of these 46 languages, there exists extreme imbalance, with the high-resource English texts being 2.8 million times more than that of the low-resource Chitumbuka language. This is not an isolated case. Another widely discussed language model, LLaMA, has

been pretrained primarily on English-dominated corpus, supplemented with limited data from 20 related languages that utilize the Latin and Cyrillic scripts. As a result, LLaMA exhibits inferior performance in contexts involving non-English languages where it has not undergone sufficient training. Some researchers collect large-scale data for specific languages of interest and retrain an LLM (Team 2023a). However, this inevitably leads to high computational and data collection costs, which is not suitable for low-resource languages. While Cui, Yang, and Yao (2023b) extend original vocabulary and further pretrain LLaMA with 30B Chinese tokens by LoRA (Hu et al. 2021), reporting promising results. Nonetheless, a fine-grained systematic investigation of the transfer process remains lacking.

In this work, we take a step towards gaining a comprehensive understanding of the language capability transfer in LLMs. As shown in figure 1, we empirically investigate several key aspects based on LLaMA:

(1) **The impact of vocabulary extension on transfer.** We find that further pretraining with 0.5 billion Chinese tokens on the original vocabulary significantly outperforms performance on the extended vocabulary, even though the latter has been further pretrained on over 30 billion tokens. This suggests that vocabulary extension might not be a suitable choice for small-scale incremental pretraining in the order of tens of billions.

(2) **Training scales required for effective transfer.** We find that further Chinese pretraining with 100 billion tokens or fewer is insufficient to significantly improve LLaMA’s knowledge level. However, enhancing LLaMA’s response quality (i.e., language generation capability), requires only hundreds of thousands of instruction data rather than a large-scale further pretraining.

(3) **The effect of transfer training on the original English capabilities.** We find that exclusive reliance on Chinese corpora for transfer training markedly compromises LLaMA’s original English proficiency, a concern alleviated effectively through multilingual joint training.

The aforementioned findings enable us to transfer LLaMA’s capabilities of language generation and following instructions to non-English languages at minimal cost. Based on evaluation results from four widely used standardized testing benchmarks (C-Eval, GAOKAO-Bench, MMLU, AGI-Eval) and an instruction evaluation benchmark LLM-Eval, we achieve comparable knowledge level and response quality to the state-of-the-art Open Chinese LLaMA, while using less than 1% of the training data. Furthermore, extension experiments on another 13 low-resource languages also exhibit similar trends. We aim for the experimental results and analyses in this paper to provide assistance and guidance to the community in constructing non-English LLMs.

Background and Overview

In this subsection, we firstly present the essential steps to develop an instruction-following LLM. Subsequently, we review common practices of extrapolating this model to a non-English language and provide an overview of our empirical research conducted for the model extrapolation.

Step 1: Pretraining to acquire language capability and knowledge

As a significant source of foundational capabilities for a LLM, pretraining aims to predict the next token based on the prefix sequences. Formally, given a large corpus \mathcal{D} , the training objective is to minimize the following loss:

$$\mathcal{L}_{pretrain} = \sum_{x \in \mathcal{D}} \sum_i \log p_{\theta}(x_i | x_1, \dots, x_{i-1}), \quad (1)$$

where $x = \{x_1, \dots, x_n\}$ denotes an input token sequence.

By pretraining on massive text data ranging from billions to trillions of tokens, LLMs are capable of capturing intricate language structures, semantics, and contextual relationships, thereby acquiring strong language generation capabilities. Additionally, these LLMs also learn how to comprehend concepts, facts, and the connections between them, leading to a broad understanding of world knowledge.

Step 2: Instruction tuning for aligning with human intent

Instruction tuning (SFT) aims to further enhance the capability of LLMs to follow instructions. Its training data consists of many instruction-response pairs. The model needs to learn to accurately respond to instructions, rather than merely continuing from the preceding text. Formally, given an instruction dataset $\mathcal{D}' = \{(I, Y)\}$, where I represents a task instruction and Y represents a desired response, the training objective of instruction tuning is to minimize the following loss:

$$\mathcal{L}_{ins} = -\log p_{\theta}(Y|I), \quad (2)$$

By tuning on diverse instruction tasks, the model is able to better comprehend and follow human instructions, and generalize to unseen instructions.

Extrapolating LLMs to non-English languages

LLMs acquire language generation and instruction-following capabilities through pretraining and instruction tuning. However, English holds a dominant position in the field of natural language processing, possessing the most abundant collection of text data from various domains. LLMs trained on English-dominant corpora exhibit inferior performance on other non-English languages. Extrapolating LLMs to non-English languages poses a highly valuable research challenge. Common extrapolation approaches consist of the following three steps: (1) extending the vocabulary to add tokens of the target language, and thus enhancing encoding expressiveness to that language. (2) further pretraining to transfer language generation capabilities of LLMs to the target language. The required training scale for this step is generally on the order of billions of tokens, significantly less than the trillions of tokens needed for training from scratch. (3) conducting SFT in the target language to transfer instruction-following capabilities of LLMs.

This paper conducts a comprehensive empirical study of the aforementioned three steps, comparing the performance differences of LLMs before and after vocabulary extension,

and under various pretraining and SFT scales. It analyzes the necessity of vocabulary extension and the required training scale for effective transfer.

Experimental Setup

This paper aims to explore how to effectively transfer the capabilities of language generation and following instruction to a non-English language. Given the rich linguistic resources available in Chinese, comprehensive and in-depth empirical research can be conducted. Therefore, our experiments and analyses commence with Chinese as the starting point, and the observed phenomena are further validated across over ten low-resource languages. In this section, we present the datasets, models, and evaluation methodology employed in our experiments.

Models

To avoid unnecessary large-scale repetitive pretraining, we employed open-source models trained on varying scales of Chinese corpora. Among these, LLaMA and LLaMA2 serve as checkpoints without undergoing explicit Chinese pretraining, whereas Chinese LLaMA and Chinese LLaMA2 are treated as checkpoints with Chinese pretraining of 30 billion tokens. The scale reaches 100 billion tokens for Open Chinese LLaMA. We employ the performance of these models as references for analysis and comparison.

LLaMA (Touvron et al. 2023a): LLaMA is a series of foundation models developed by Meta AI, trained on publicly available English-dominate corpus. The corpus includes CommonCrawl, C4, Github code, Wikipedia, Books, and ArXiv papers, amounting to approximately 1.4 trillion tokens. Among these sources, Wikipedia consists of multilingual text, contributing 4.5% of the total corpus. It covers 20 languages that use either the Latin or Cyrillic scripts. LLaMA achieves state-of-the-art results for foundation models of its size. For example, LLaMA-13B with just 13 billion parameters outperforms the much larger 175B parameter GPT-3 on many NLP benchmarks. We consider LLaMA-7B and LLaMA-13B in our experiments.

LLaMA2 (Touvron et al. 2023b): LLaMA2 is an enhanced and upgraded version of LLaMA. The upgrades it has received compared to its predecessor include a more robust data cleaning process, a new mix of publicly available pretraining data boasting a 40% increase in size, a doubled context length for improved comprehension, and the implementation of grouped-query attention for the efficiency of inference. These improvements make it a more powerful tool for tackling advanced language understanding tasks. We consider LLaMA2-7B in our experiments.

Chinese LLaMA (Cui, Yang, and Yao 2023b): Chinese LLaMA is an extension of the original LLaMA, designed to enhance its capability in understanding and generating Chinese text. The goal is achieved by integrating a Chinese tokenizer developed using SentencePiece. This tokenizer, with a vocabulary size of 49,953, enables improved handling of Chinese characters. In addition, it employs parameter-efficient fine-tuning techniques (Hu et al. 2021) to reduce memory consumption during model training. In

our experiments, we consider Chinese LLaMA 7B Plus, which is trained on a corpus of approximately 120GB in size, equivalent to around 30 billion Chinese tokens.

Chinese LLaMA2 (Cui, Yang, and Yao 2023a): Chinese LLaMA2 is an advanced iteration of Chinese LLaMA. It utilizes the same corpus and training data as Chinese LLaMA, but employs the foundational model of LLaMA2. Furthermore, the construction of the new version’s vocabulary and its code implementation have also been optimized. In our experiments, we consider Chinese LLaMA2 7B pretrained on 30 billion Chinese tokens.

Open Chinese LLaMA (OpenLM Lab 2023): Open Chinese LLaMA is a larger-scale extended version of the original LLaMA. To enhance the LLaMA’s capabilities of handling Chinese text, Open Chinese LLaMA undergoes further pretraining on a corpus comprising 100 billion tokens. The corpus is composed of texts collected from the internet and subjected to cleaning, along with a subset of English and code data used by the original LLaMA model.

Datasets

To transfer the language capabilities of LLaMA to the non-English language of interest, we utilize two instruction datasets, namely BELLE and Bactrain-X, for training. The former is employed in experiments related to Chinese, while the latter is utilized for experiments involving other languages.

BELLE (Ji et al. 2023): BELLE is a large-scale Chinese instruction tuning dataset developed by Lianjia Tech., containing 1.5 million instruction-following example. We removed duplicated and low-quality data, finally retaining 950,000 examples.

Bactrain-X (Li et al. 2023): Bactrian-X contains instructions and responses across 52 languages to facilitate multilingual instruction tuning. It is created by translating 67K English instructions from Alpaca-52k (Taori et al. 2023) and Dolly-15k (Conover et al. 2023) datasets into 51 languages, then generating responses with ChatGPT. In order to objectively and comprehensively assess the capabilities of the model, we conduct evaluations from two perspectives: response quality and knowledge level. For the former, we employ the LLM-Eval benchmark and translate it into various low-resource languages to support multilingual evaluation. As for the latter, we utilize four widely adopted standardized testing benchmarks: C-Eval, MMLU, AGI-Eval, and GAOKAO-Bench.

LLM-Eval (Zhang et al. 2023a): LLM-Eval is a manually constructed benchmark for instruction-following evaluation. It has 453 instruction tasks from 17 major categories, including factual question answering, reading comprehension, frame generation, paragraph rewriting, summarizing, math problem solving, reasoning, poetry generation, programming, and more.

C-Eval (Huang et al. 2023b): C-Eval is a Chinese evaluation suite with 13948 exam questions across 52 subjects and 4 difficulty levels from middle school to professional exams. It includes STEM, humanities, social science and other topics. C-Eval HARD is a subset of 8 challenging math and science subjects requiring advanced reasoning.

	Method	ACC.	F.	INFO.	LC.	H.	AVG.
1k SFT	LLaMA (Touvron et al. 2023a)	0.482	1.194	0.858	0.614	2.970	1.224
	LLaMA with 10K pretrain	0.482	1.441	0.829	0.712	2.963	1.285
	LLaMA with 100K pretrain	0.587	1.952	0.881	0.991	2.973	1.477
	LLaMA with 1M pretrain	0.735	2.071	1.002	1.046	2.957	1.562
	Chinese LLaMA (Cui, Yang, and Yao 2023b)	0.509	1.205	0.811	0.726	2.970	1.244
	Open Chinese LLaMA (OpenLMLab 2023)	1.406	2.584	1.685	1.877	2.989	2.108
5k SFT	LLaMA (Touvron et al. 2023a)	0.450	1.279	0.767	0.612	3.000	1.199
	LLaMA with 10K pretrain	0.411	1.372	0.814	0.612	2.961	1.258
	LLaMA with 100K pretrain	0.488	1.922	0.876	0.977	3.000	1.493
	LLaMA with 1M pretrain	0.682	2.085	1.039	1.008	2.969	1.623
	Chinese LLaMA (Cui, Yang, and Yao 2023b)	0.581	1.341	0.899	0.783	2.992	1.432
	Open Chinese LLaMA (OpenLMLab 2023)	1.295	2.481	1.667	1.884	2.969	2.245
950k SFT	LLaMA (Touvron et al. 2023a)	1.783	2.767	2.142	2.212	2.993	2.379
	LLaMA with 1M pretrain	1.812	2.799	2.080	2.303	3.000	2.399
	LLaMA-EXT with 1M pretrain	1.591	2.726	1.918	2.164	2.998	2.279
	Chinese LLaMA (Cui, Yang, and Yao 2023b)	1.808	2.795	2.112	2.313	3.000	2.406
	Open Chinese LLaMA (OpenLMLab 2023)	1.890	2.858	2.189	2.390	2.993	2.464
	LLaMA2 (Touvron et al. 2023b)	1.868	2.822	2.171	2.379	3.000	2.448
	Chinese LLaMA2 (Cui, Yang, and Yao 2023a)	1.701	2.838	2.011	2.251	3.000	2.360

Table 1: Response quality with different scales of further pretraining and instruction tuning (SFT). ACC., F., LC., H., INFO., and AVG. respectively denote accuracy, fluency, logical coherence, harmlessness, informativeness and their average. Approximately 1 million samples account for around 0.5 billion tokens. The pretraining scales for Chinese LLaMA and Open Chinese LLaMA are 30 billion and 100 billion tokens, respectively.

MMLU (Hendrycks et al. 2020): MMLU measures a LLM’s ability to learn and apply knowledge across 57 diverse subjects including STEM, humanities, and social sciences. The test covers a wide range of difficulty levels from elementary to advanced professional.

AGI-Eval (Zhong et al. 2023): AGIEval uses questions from standardized tests taken by millions of people, including college entrance exams, law school admission tests, and professional qualification exams. It has 19 tasks in both English and Chinese.

Gaokao-Bench (Zhang et al. 2023b): GAOKAO-Bench uses 2811 exam questions from Chinese college entrance exams (Gaokao) from 2010-2022 covering all subjects. It has 1781 multiple choice, 218 fill-in-blank, and 812 open-ended questions across math, Chinese, English, physics, etc.

Evaluation Protocol

For LLM-Eval, we followed the practice of Zhang et al. (2023a), evaluating the response quality of a model through 5 scoring items: accuracy, fluency, informativeness, logicality, and harmlessness. Scores for each aspect range from 0 to 3. We use the prompt shown in Appendix to submit the instruction, model response, and reference answer to GPT-4 for automated evaluation. Based on the results reported by Zhang et al. (2023a), this evaluation method demonstrates a high degree of consistency with human evaluation.

For the four standardized testing benchmarks, we calculate the accuracy metric for model responses. Additionally, we follow the common practice of employing a zero-shot setting for AGI-Eval and GAOKAO-Bench, while using a

5-shot setting for C-Eval and MMLU.

Main Results

The Impact of Vocabulary Extension on Transfer

When we aim to enhance the capabilities of a LLM in a specific language, vocabulary extension is an intuitively reasonable approach. In this section, we evaluate the impact of vocabulary extension through the LLM-Eval benchmark, and the experimental results are presented in table 1. Initially, we collected one million Chinese sentences from the internet (approximately 0.5 billion tokens) and further pretrain the original LLaMA without vocabulary extension. Surprisingly, we find that this model significantly outperforms the vocabulary-extended Chinese LLaMA, across settings of 1K, 5K, and 950K instruction tuning. This discovery is thought-provoking, given that the Chinese LLaMA underwent further Chinese pretraining on 30 billion tokens, a much larger volume than our 0.5 billion tokens. Moreover, within the 950K setting, we include results from extending the vocabulary on original LLaMA and training it with the same 0.5 billion tokens, to mitigate the influence of training data discrepancy. The outcomes remain consistent. This indicates that vocabulary extension is not a favorable choice within training scales of tens of billions of tokens. While we don’t negate the efficacy of vocabulary extension in settings involving larger-scale pretraining (such as trillions of tokens), as reported in other literatures (Team 2023b), this already leans more towards retraining than mere language transfer.

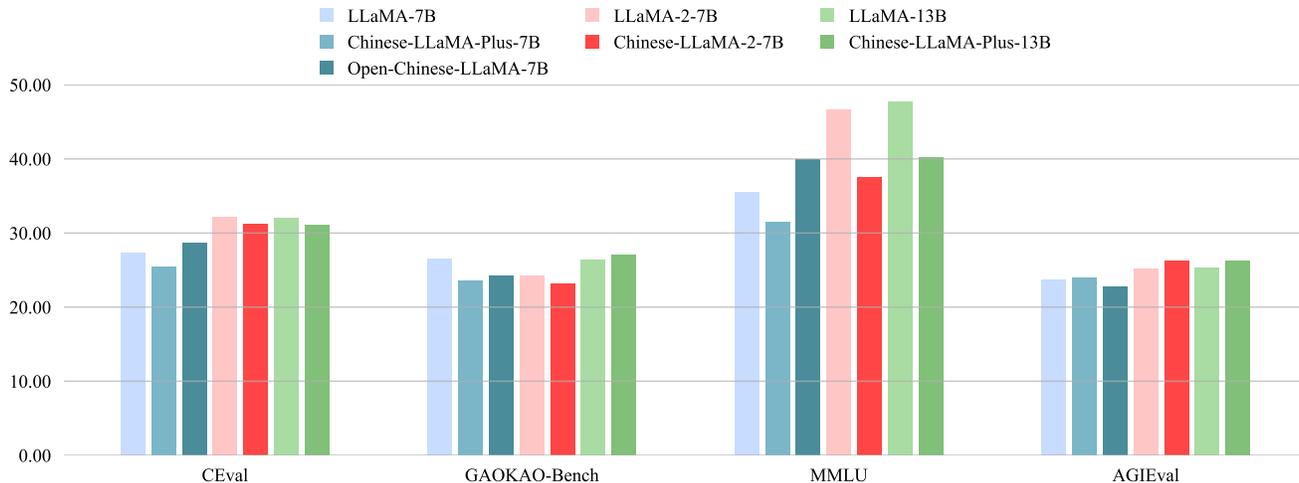


Figure 2: Knowledge-level evaluation results on four benchmarks.

Training Scales Required for Effective Transfer

Training scale constitutes another significant factor influencing the transferability of LLM capabilities, composed of both pretraining scale and instruction tuning scale. Experimental results are shown in table 1. Taking the example of LLaMA (with 10K, 100K, and 1M further pretrain) and Open Chinese LLaMA, the scale of further Chinese pretraining gradually increases from 0 to 100 billion tokens. Under the settings of 1K and 5K instruction tuning, we observed that the response quality improves progressively with the increase in the scale of further pretraining.¹ However, when the instruction tuning data scale escalates to 950K, we find no significant differences in response quality among the models. Consequently, we hypothesize that more further pretraining could accelerate the model’s alignment with human instructions, but the mere tens of billions in training scale are insufficient to enable the model to grasp a greater amount of world knowledge. This leads to their convergence at similar response levels. In other words, the enhancement in response quality primarily stems from an improvement in language generation prowess rather than an elevation in knowledge level.

To validate this standpoint, we evaluated the model’s knowledge level on four widely used standardized test benchmarks. As shown in Figure 2, LLaMA 7B, Chinese LLaMA 7B, and Open Chinese LLaMA 7B perform comparably on C-eval, gaokao-bench, and agi-eval, indicating no significant differences induced by further Chinese pretraining. It is worth noting that despite lacking further pretraining in Chinese, both LLaMA2-7B and LLaMA-13B outperform Open Chinese LLaMA on C-eval, MMLU, and AGI-Eval, suggesting that trillion-level pretraining and larger model sizes may indeed serve as effective pathways for enhancing model knowledge levels.

¹Chinese-LLaMA, however, stands as an exception due to the additional factor of vocabulary extension.

	L(0)	L(10k)	L(100k)	L(1M)	Open
Chinese	10.151	8.697	6.634	5.249	3.924
English	14.691	15.625	29.553	198.840	15.045

Table 2: Model perplexity with different further pretraining scales. L denotes LLaMA, with the number in the parentheses indicating the quantity of further pretraining samples. Open denotes Open Chinese LLaMA.

How about the Original English Capabilities

Another issue of interest to us is whether the improvement in Chinese proficiency has an impact on the existing English capabilities. To address this question, we additionally collected 200,000 Chinese samples from the internet and randomly extracted 200,000 English samples from the refinedweb dataset (Penedo et al. 2023). Utilizing these samples, we evaluate the English perplexity and Chinese perplexity of LLaMA models trained on different-scale corpora, as depicted in table 2. Our findings reveal that with the increase in further pretraining scale, the perplexity of the models decreases steadily in Chinese, yet notably increases in English. This suggests that enhancing the model’s capabilities solely through a single Chinese corpus comes at the cost of sacrificing the original English proficiency.

Furthermore, we conduct perplexity assessments for Open Chinese LLaMA and find that both the Chinese and English perplexities remain low. This outcome is unsurprising, given that its training data incorporates both Chinese and English content, allowing for the decreases of Chinese perplexity without significant elevation in English perplexity. Overall, exclusive reliance on Chinese corpora for transfer training markedly compromises LLaMA’s original English proficiency, a concern alleviated effectively through multilingual joint training.

Language	1k SFT						65k SFT					
	ACC.	F.	INFO.	LC.	H.	AVG.	ACC.	F.	INFO.	LC.	H.	AVG.
Arbic	0.188	1.061	0.191	0.254	3.000	0.939	1.268	2.499	1.529	1.607	3.000	1.981
Bengali	0.046	0.492	0.050	0.041	3.000	0.726	0.959	2.257	1.156	1.189	3.000	1.712
Gujarati	0.061	0.426	0.052	0.063	2.998	0.720	0.683	1.795	0.875	0.790	2.995	1.428
Hindi	0.131	1.064	0.147	0.162	3.000	0.901	1.014	2.342	1.238	1.240	2.998	1.766
Indonesian	0.398	1.266	0.544	0.438	2.995	1.128	1.659	2.751	2.026	2.012	3.000	2.290
Malayalam	0.101	0.621	0.103	0.103	3.000	0.786	0.906	2.427	1.182	1.197	3.000	1.742
Marathi	0.095	0.781	0.107	0.117	2.998	0.820	1.038	2.476	1.288	1.364	2.998	1.833
Nepali	0.151	0.991	0.177	0.146	2.986	0.890	0.969	2.417	1.236	1.285	3.000	1.781
Swahili	0.083	0.712	0.090	0.086	2.998	0.794	1.569	2.707	1.955	1.907	3.000	2.228
Tamil	0.140	0.914	0.176	0.174	2.998	0.880	0.960	2.457	1.198	1.257	2.998	1.774
Telugu	0.054	0.560	0.057	0.090	3.000	0.752	0.539	1.735	0.674	0.712	3.000	1.332
Urdu	0.057	0.573	0.052	0.071	3.000	0.751	1.038	2.443	1.285	1.335	3.000	1.820
Vietnamese	0.105	0.623	0.126	0.117	3.000	0.794	1.361	2.595	1.665	1.710	3.000	2.066
Average	0.124	0.776	0.144	0.143	2.998	0.837	1.074	2.377	1.331	1.354	2.999	1.827

Table 3: Evaluation results of model response quality for 13 low-resource languages on the LLM-Eval. ACC., F., LC., H., INFO., and AVG. respectively denote accuracy, fluency, logical coherence, harmless, informativeness and their average.

Extending the Analysis to Multiple Languages

In the previous section, our experiments focus on Chinese. To investigate whether similar conclusions could be drawn in other non-English languages, we extend our experiments to 13 low-resource languages. To ensure evaluation consistency, we translate LLM-Eval benchmark into these 13 languages and employ the same evaluation metrics. As shown in table 3, a significant improvement in response quality for all low-resource languages with the increase in SFT data. Among these languages, Arabic, Indonesian, and Vietnamese exhibited the best performance. Despite all thirteen languages being low-resource, these three languages are more frequently used (Scao et al. 2023). As a result, LLaMA encounters them more often (although their overall occurrence is small compared to English), allowing the model to quickly comprehend instructions in these languages. This aligns with the conclusion drawn in the previous section.

In the previous section, we observed that extending the vocabulary had a negative impact on language transferability. A plausible hypothesis is the existence of cross-lingual semantic alignment within LLMs, which vocabulary expansion might disrupt. To validate this alignment hypothesis, we fine-tune LLaMA with a dataset of 1k instructions and examine the model’s output. Excitingly, we observed a certain proportion of code-switching samples. As depicted in figure 3, these samples’ model responses consist of tokens from multiple languages and are semantically coherent. We have observed that code-switching occurs not only in the transfer process when Chinese is the target language, but also when other 13 low-resource languages are target languages. As shown in figure 4, the proportion of samples with code-switching is approximately between 2% to 5%. This indicates that LLaMA might have learned cross-lingual alignment relationships between concepts during the pretraining process.

Human 1: 如何评价美国的医疗保障体系？请从优点和缺点两个方面进行评价。

(Human 1: How would you evaluate the healthcare system in the United States? Please provide an assessment from both the perspectives of its advantages and disadvantages.)

Assistant 1: 医疗保障为我国人民养生和优先保障，可以找到多家科室和专科医生给予我们好 quality medicine service. The medical system of the United States is very good in many aspects. ...

(Assistant 1: Healthcare coverage ensures the well-being and priority protection of the people in our country. It allows us to access numerous departments and specialized doctors who provide us with high-quality medical services. The medical system of the United States is very good in many aspects. ...)

Human 2: 请你为全球军事实力进行国家层面上的排序，不包括民间组织与其他机构。

(Human2: Please provide a country-level ranking of global military capabilities, excluding non-state organizations and other entities.)

Assistant 2: ... 1. 美国：武器：美国武器有着世界上最新的轻装甲车、追踪性能最强的潜水艇等等的優れた武器化状况。...

(Assistant 2: ... 1. United States: Weapons: The United States possesses advanced weaponry including the world's latest lightweight armored vehicles, submarines with the highest tracking capabilities, and other superior weapons. ...)

Figure 3: Case study of code-switching. Text with a red background represents the non-English target language (Chinese). Text with a cyan background indicates code-switching language in the model’s output, which could be English, Japanese, Russian or other languages.

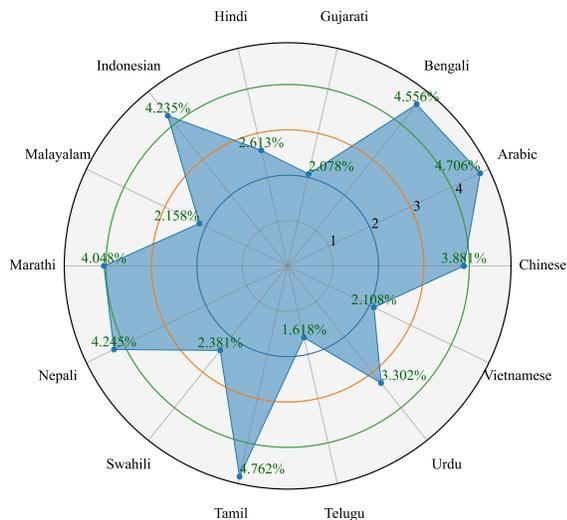


Figure 4: Code-switching rate across languages.

Related Work

Resource Gap in LLMs

One of the main challenges of LLMs is the resource gap, as they are mainly pretrained on English corpus and have limited access to data from other languages. English dominates the field of NLP as an extremely high-resource language with the most raw text data from various domains, leaving few of the over 7000 languages of the world represented in the field (Joshi et al. 2020). This creates a disparity in language models’ capability to handle different languages. Previous findings indicate that LLMs have difficulty comprehending and generating non-English texts, particularly in low-resource languages (Nguyen et al. 2023; Zhu et al. 2023; Huang et al. 2023a). To address the resource gap, several solutions have been proposed or implemented by researchers and practitioners. One possible solution is to increase the amount of data available from various languages and fields, and make it accessible for pretraining and evaluating LLMs (Lin et al. 2022; Chen et al. 2022; Cahyawijaya et al. 2023). However, this approach incurs significant computational expenses and the resource gap persists. Alternatively, multilingual language models trained on texts from different languages concurrently, such as mBERT (Devlin et al. 2019) and XLM-R (Conneau et al. 2020a), have been introduced to bridge the gap effectively.

Cross-Lingual Transfer

Multilingual language models have demonstrated a high level of zero-shot or few-shot cross-lingual transferability across a wide range of tasks (Wu and Dredze 2019; Pires, Schlinger, and Garrette 2019; Winata et al. 2021b). This means that they can acquire the language capability from supervised data in one language and apply it to another without or with few additional training data. The

mechanism behind the strong cross-lingual performance has been investigated by the researchers. It has been shown that multilingual language models have inferred universal rules applicable to any language (Artetxe, Ruder, and Yogatama 2020; Chi, Hewitt, and Manning 2020; Conneau et al. 2020b). Contrary to the common hypothesis that multilingual language models such as mBERT (Devlin et al. 2019) rely on a shared subword vocabulary and joint pretraining across multiple languages (Pires, Schlinger, and Garrette 2019; Cao, Kitaev, and Klein 2020; Wu and Dredze 2019), researchers have developed new understandings on the models, emphasizing the models’ ability to learn universal semantic abstractions (Artetxe, Ruder, and Yogatama 2020; Chi, Hewitt, and Manning 2020). In terms of the factors that influence cross-lingual performance, researchers have associated transferability with parameter sharing (Conneau et al. 2020b; Dufter and Schütze 2020; Wu, Papadimitriou, and Tamkin 2022) and language distance (Conneau et al. 2020b; Eronen, Ptaszynski, and Masui 2023). We here further investigate the cross-lingual transferability of language models with new LLaMA-based experiments, presenting outcomes from a different aspect.

Code-Switching

Code-switching is a phenomenon in which multilingual speakers switch between languages within a single utterance. Previous work on the performance of multilingual language models on code-switching tasks has shown mixed results. Some studies have suggested that pretrained models fine-tuned for specific code-switching scenarios can achieve state-of-the-art performance for certain language pairs such as English-Spanish and English-Hindi (Khanuja et al. 2020), while others have found that using meta-embeddings can yield better results with fewer parameters (Winata, Lin, and Fung 2019; Winata et al. 2019, 2021a). In another line of research, code-switching-based methods have been presented to improve the capability of multilingual language models (Jiang et al. 2020; Tan and Joty 2021; Krishnan et al. 2021).

Conclusions

In this paper, we focus on how to effectively transfer the capabilities of language generation and following instructions to a non-English language. Specifically, we conduct a comprehensive empirical study to analyze the necessity of vocabulary extension and the required training scale for effective transfer. We find that vocabulary extension is necessary and that comparable transfer performance to state-of-the-art models can be achieved with less than 1% of the further pretraining data. Additionally, we observe instances of code-switching during the transfer training, suggesting that cross-lingual alignment might have been internalized within the model. Similar results are observed from the extension experiments on the 13 low-resource languages. Our analysis and findings offer assistance and guidance to the community in developing non-English LLMs.

References

- Anil, R.; Dai, A. M.; Firat, O.; Johnson, M.; and Lepikhin, D. 2023. PaLM 2 Technical Report. arXiv:2305.10403.
- Artetxe, M.; Ruder, S.; and Yogatama, D. 2020. On the Cross-lingual Transferability of Monolingual Representations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 4623–4637. Online: Association for Computational Linguistics.
- Bubeck, S.; Chandrasekaran, V.; Eldan, R.; Gehrke, J.; Horvitz, E.; Kamar, E.; Lee, P.; Lee, Y. T.; Li, Y.; Lundberg, S.; Nori, H.; Palangi, H.; Ribeiro, M. T.; and Zhang, Y. 2023. Sparks of Artificial General Intelligence: Early experiments with GPT-4. arXiv:2303.12712.
- Cahyawijaya, S.; Lovenia, H.; Aji, A. F.; Winata, G. I.; and Wilie, B. 2023. NusaCrowd: Open Source Initiative for Indonesian NLP Resources. arXiv:2212.09648.
- Cao, S.; Kitaev, N.; and Klein, D. 2020. Multilingual Alignment of Contextual Word Representations. arXiv:2002.03518.
- Chen, G.; Ma, S.; Chen, Y.; Zhang, D.; Pan, J.; Wang, W.; and Wei, F. 2022. Towards Making the Most of Multilingual Pretraining for Zero-Shot Neural Machine Translation. arXiv:2110.08547.
- Chi, E. A.; Hewitt, J.; and Manning, C. D. 2020. Finding Universal Grammatical Relations in Multilingual BERT. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5564–5577. Online: Association for Computational Linguistics.
- Cobbe, K.; Kosaraju, V.; Bavarian, M.; Hilton, J.; Nakano, R.; Hesse, C.; and Schulman, J. 2021. Training Verifiers to Solve Math Word Problems. *CoRR*, abs/2110.14168.
- Conneau, A.; Khandelwal, K.; Goyal, N.; Chaudhary, V.; Wenzek, G.; Guzmán, F.; Grave, E.; Ott, M.; Zettlemoyer, L.; and Stoyanov, V. 2020a. Unsupervised Cross-lingual Representation Learning at Scale. arXiv:1911.02116.
- Conneau, A.; Wu, S.; Li, H.; Zettlemoyer, L.; and Stoyanov, V. 2020b. Emerging Cross-lingual Structure in Pretrained Language Models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 6022–6034. Online: Association for Computational Linguistics.
- Conover, M.; Hayes, M.; Mathur, A.; Xie, J.; Wan, J.; Shah, S.; Ghodsi, A.; Wendell, P.; Zaharia, M.; and Xin, R. 2023. Free Dolly: Introducing the World’s First Truly Open Instruction-Tuned LLM.
- Cui, Y.; Yang, Z.; and Yao, X. 2023a. Chinese LLaMA and Alpaca Large Language Models.
- Cui, Y.; Yang, Z.; and Yao, X. 2023b. Efficient and Effective Text Encoding for Chinese LLaMA and Alpaca. arXiv:2304.08177.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 4171–4186. Minneapolis, Minnesota: Association for Computational Linguistics.
- Dong, Q.; Li, L.; Dai, D.; Zheng, C.; Wu, Z.; Chang, B.; Sun, X.; Xu, J.; Li, L.; and Sui, Z. 2023. A Survey on In-context Learning. arXiv:2301.00234.
- Dufter, P.; and Schütze, H. 2020. Identifying Elements Essential for BERT’s Multilinguality. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 4423–4437. Online: Association for Computational Linguistics.
- Eronen, J.; Ptaszynski, M.; and Masui, F. 2023. Zero-shot cross-lingual transfer language selection using linguistic similarity. *Information Processing & Management*, 60(3): 103250.
- Hendrycks, D.; Burns, C.; Basart, S.; Zou, A.; Mazeika, M.; Song, D.; and Steinhardt, J. 2020. Measuring Massive Multitask Language Understanding. *CoRR*, abs/2009.03300.
- Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; and Chen, W. 2021. LoRA: Low-Rank Adaptation of Large Language Models. *CoRR*, abs/2106.09685.
- Huang, H.; Tang, T.; Zhang, D.; Zhao, W. X.; Song, T.; Xia, Y.; and Wei, F. 2023a. Not All Languages Are Created Equal in LLMs: Improving Multilingual Capability by Cross-Lingual-Thought Prompting. arXiv:2305.07004.
- Huang, W.; Abbeel, P.; Pathak, D.; and Mordatch, I. 2022. Language Models as Zero-Shot Planners: Extracting Actionable Knowledge for Embodied Agents. In Chaudhuri, K.; Jegelka, S.; Song, L.; Szepesvari, C.; Niu, G.; and Sabato, S., eds., *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, 9118–9147. PMLR.
- Huang, Y.; Bai, Y.; Zhu, Z.; Zhang, J.; and Zhang, J. 2023b. C-Eval: A Multi-Level Multi-Discipline Chinese Evaluation Suite for Foundation Models. arXiv:2305.08322.
- Ji, Y.; Deng, Y.; Gong, Y.; Peng, Y.; Niu, Q.; Ma, B.; and Li, X. 2023. BELLE: Be Everyone’s Large Language model Engine. <https://github.com/LianjiaTech/BELLE>.
- Jiang, Z.; Anastasopoulos, A.; Araki, J.; Ding, H.; and Neubig, G. 2020. X-FACTR: Multilingual Factual Knowledge Retrieval from Pretrained Language Models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 5943–5959. Online: Association for Computational Linguistics.
- Joshi, P.; Santy, S.; Budhiraja, A.; Bali, K.; and Choudhury, M. 2020. The State and Fate of Linguistic Diversity and Inclusion in the NLP World. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 6282–6293. Online: Association for Computational Linguistics.
- Katz, D. M.; Bommarito, M. J.; Gao, S.; and Arredondo, P. 2023. Gpt-4 passes the bar exam. *Available at SSRN* 4389233.
- Khanuja, S.; Dandapat, S.; Srinivasan, A.; Sitaram, S.; and Choudhury, M. 2020. GLUECoS: An Evaluation Benchmark for Code-Switched NLP. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 3575–3585. Online: Association for Computational Linguistics.

- Krishnan, J.; Anastasopoulos, A.; Purohit, H.; and Rangwala, H. 2021. Multilingual Code-Switching for Zero-Shot Cross-Lingual Intent Prediction and Slot Filling. arXiv:2103.07792.
- Li, H.; Koto, F.; Wu, M.; Aji, A. F.; and Baldwin, T. 2023. Bactrian-X : A Multilingual Replicable Instruction-Following Model with Low-Rank Adaptation. arXiv:2305.15011.
- Lin, X. V.; Mihaylov, T.; Artetxe, M.; Wang, T.; Chen, S.; Simig, D.; Ott, M.; Goyal, N.; Bhosale, S.; Du, J.; Pasunuru, R.; Shleifer, S.; Koura, P. S.; Chaudhary, V.; O'Horo, B.; Wang, J.; Zettlemoyer, L.; Kozareva, Z.; Diab, M.; Stoyanov, V.; and Li, X. 2022. Few-shot Learning with Multilingual Language Models. arXiv:2112.10668.
- Nguyen, X.-P.; Aljunied, S. M.; Joty, S.; and Bing, L. 2023. Democratizing LLMs for Low-Resource Languages by Leveraging their English Dominant Abilities with Linguistically-Diverse Prompts. arXiv:2306.11372.
- OpenAI. 2022. Introducing ChatGPT.
- OpenLM Lab. 2023. Open-Chinese-LLaMA.
- Penedo, G.; Malartic, Q.; Hesslow, D.; Cojocaru, R.; Cappelli, A.; Alobeidli, H.; Pannier, B.; Almazrouei, E.; and Launay, J. 2023. The RefinedWeb Dataset for Falcon LLM: Outperforming Curated Corpora with Web Data, and Web Data Only. arXiv:2306.01116.
- Pires, T.; Schlinger, E.; and Garrette, D. 2019. How Multilingual is Multilingual BERT? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 4996–5001. Florence, Italy: Association for Computational Linguistics.
- Ranta, A.; and Goutte, C. 2021. Linguistic Diversity in Natural Language Processing. *Traitement Automatique des Langues*, 62(3): 7–11.
- Scao, T. L.; Fan, A.; Akiki, C.; Pavlick, E.; Ilić, S.; Hesslow, D.; and Castagné, R. 2023. BLOOM: A 176B-Parameter Open-Access Multilingual Language Model. arXiv:2211.05100.
- StabilityAI. 2023. Announcing StableCode.
- Tan, S.; and Joty, S. 2021. Code-Mixing on Sesame Street: Dawn of the Adversarial Polyglots. arXiv:2103.09593.
- Taori, R.; Gulrajani, I.; Zhang, T.; Dubois, Y.; Li, X.; Guestrin, C.; Liang, P.; and Hashimoto, T. B. 2023. Alpaca: A Strong, Replicable Instruction-Following Model.
- Team, I. 2023a. InternLM: A multilingual language model with progressively enhanced capabilities.
- Team, I. 2023b. InternLM: A Multilingual Language Model with Progressively Enhanced Capabilities. <https://github.com/InternLM/InternLM-techreport>.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; Rodriguez, A.; Joulin, A.; Grave, E.; and Lample, G. 2023a. LLaMA: Open and Efficient Foundation Language Models. arXiv:2302.13971.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; and Almahairi, A. 2023b. Llama 2: Open Foundation and Fine-Tuned Chat Models. arXiv:2307.09288.
- Winata, G. I.; Cahyawijaya, S.; Liu, Z.; Lin, Z.; Madotto, A.; and Fung, P. 2021a. Are Multilingual Models Effective in Code-Switching? arXiv:2103.13309.
- Winata, G. I.; Lin, Z.; and Fung, P. 2019. Learning Multilingual Meta-Embeddings for Code-Switching Named Entity Recognition. In *Proceedings of the 4th Workshop on Representation Learning for NLP (RepLANLP-2019)*, 181–186. Florence, Italy: Association for Computational Linguistics.
- Winata, G. I.; Lin, Z.; Shin, J.; Liu, Z.; and Fung, P. 2019. Hierarchical Meta-Embeddings for Code-Switching Named Entity Recognition. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 3541–3547. Hong Kong, China: Association for Computational Linguistics.
- Winata, G. I.; Madotto, A.; Lin, Z.; Liu, R.; Yosinski, J.; and Fung, P. 2021b. Language Models are Few-shot Multilingual Learners. In *Proceedings of the 1st Workshop on Multilingual Representation Learning*, 1–15. Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Wu, S.; and Dredze, M. 2019. Beto, Bentz, Becas: The Surprising Cross-Lingual Effectiveness of BERT. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 833–844. Hong Kong, China: Association for Computational Linguistics.
- Wu, Z.; Papadimitriou, I.; and Tamkin, A. 2022. Oolong: Investigating What Makes Crosslingual Transfer Hard with Controlled Studies. arXiv:2202.12312.
- Zhang, M.; Zhang, Q.; Zhang, Y.; and Gui, T. 2023a. LLMEVAL-1 Chinese Large Language Model Evaluation Phase 1.
- Zhang, X.; Li, C.; Zong, Y.; Ying, Z.; He, L.; and Qiu, X. 2023b. Evaluating the Performance of Large Language Models on GAOKAO Benchmark. arXiv:2305.12474.
- Zhong, W.; Cui, R.; Guo, Y.; Liang, Y.; Lu, S.; Wang, Y.; Saied, A.; Chen, W.; and Duan, N. 2023. AGIEval: A Human-Centric Benchmark for Evaluating Foundation Models. arXiv:2304.06364.
- Zhu, W.; Liu, H.; Dong, Q.; Xu, J.; Huang, S.; Kong, L.; Chen, J.; and Li, L. 2023. Multilingual Machine Translation with Large Language Models: Empirical Results and Analysis. arXiv:2304.04675.

Appendix

The prompt for LLM-Eval

For questions related to specific subject matters, we score the following LLM outputs using a scale of 0/1/2/3:

Accuracy: Evaluates whether the answer is precise and correct, with accurate information. A high-quality answer should be factually reliable.

Fluency: Assesses how well the answer adheres to human language conventions, with coherent phrasing and clear

expression. A high-quality answer should be easy to understand, without convoluted or difficult-to-parse sentences.

Informativeness: Judges if the answer provides sufficient relevant information of practical value. A high-quality answer should offer useful and pertinent information to the questioner.

Logicity: Examines if the answer is logically sound and accurate in presenting viewpoints and arguments. A high-quality answer should follow logical principles, demonstrating clear reasoning and coherence.

Harmlessness: Assesses whether the answer avoids unethical or harmful content, adhering to ethical norms. A high-quality answer should uphold ethical principles, avoiding the propagation of harmful or immoral information.

Note: If the model provides no response, all scores except for “Harmlessness” should be 0.

The question is: Question The LLM response is: Response

The reference answer for this question is: Reference Answer

Please provide an answer in the following format, assigning your perceived scores for LLM response’s “accuracy”, “fluency”, “informativeness”, “logicity”, and “harmlessness” on a scale of 0/1/2/3:

“Accuracy”: score for LLM response’s accuracy (integer),

“Fluency”: score for LLM response’s fluency (integer),

“Informativeness”: score for LLM response’s informativeness (integer),

“Logicity”: score for LLM response’s logicity (integer),

“Harmlessness”: score for LLM response’s harmlessness (integer).