

# Language Versatilists vs. Specialists: An Empirical Revisiting on Multilingual Transfer Ability

Anonymous ACL submission

## Abstract

Multilingual transfer ability, which reflects how well the models fine-tuned on one source language can be applied to other languages, has been well studied in multilingual pre-trained models (e.g., BLOOM (Scao et al., 2022)). However, such ability has not been investigated for English-centric models (e.g., LLaMA (Touvron et al., 2023a)). To fill this gap, we study the following research questions. First, does multilingual transfer ability exist in English-centric models and how does it compare with multilingual pretrained models? Second, does it only appears when English is the source language for the English-centric model? Third, how does it vary in different tasks? We take multilingual reasoning ability as our focus and conduct extensive experiments across four types of reasoning tasks. We find that the multilingual pretrained model does not always outperform an English-centric model. Furthermore, English appears to be a less suitable source language, and the choice of source language becomes less important when the English-centric model scales up. In addition, different types of tasks exhibit different multilingual transfer abilities. These findings demonstrate that English-centric models not only possess multilingual transfer ability but may even surpass the transferability of multilingual pretrained models if well-trained. By showing the strength and weaknesses, the experiments also provide valuable insights into enhancing multilingual reasoning abilities for the English-centric models.

## 1 Introduction

Multilingual pre-training has become a standard technique to equip a language model with cross-lingual transfer ability, through which it is possible to improve the performance on low-resource languages by leveraging high-resource languages (Devlin et al., 2019; Conneau et al., 2018a, 2020; Lin et al., 2021; Scao et al., 2022). However, there

have been looming concerns regarding multilingual pre-training. For instance, Conneau et al. (2020) uncovered *the curse of multilinguality*, suggesting for a fixed model size, cross-lingual performance increases with additional pretraining languages only up to a certain point, after which the performance begins to decline. Additionally, Wang et al. (2020) also reported a phenomenon called *negative interference*, meaning the performance on both high-resource and low-resource languages degrade due to joint multilingual learning.

English-centric models (Brown et al., 2020; Chowdhery et al., 2022; Black et al., 2021; Wang and Komatsuzaki, 2021; Black et al., 2022; Biderman et al., 2023; Zhang et al., 2022; Touvron et al., 2023a,b), on the other hand, have demonstrated strong performance on downstream English tasks, but their cross-lingual abilities have not been systematically analyzed.<sup>1</sup> While it may seem intuitive to assume that English-centric models are not well-suited in cross-lingual transfer, this is not necessarily the case in practice. Research evidence suggests that monolingual models are capable of learning certain abstractions that can generalize across languages, as demonstrated by Artetxe et al. (2020). In addition, it should be noted that English-centric models are not limited to English only, as they have been exposed to some other languages, albeit to a much lesser extent (Brown et al., 2020; Gao et al., 2020; Chowdhery et al., 2022; Touvron et al., 2023a,b).

The investigation of multilingual models and English-centric models is especially meaningful in many practical settings. Suppose the goal is to develop a model with excellent multilingual reasoning skills such as arithmetic, commonsense, and logical reasoning. In that case, how should we approach this goal? Should we start from

<sup>1</sup>In this paper, we refer to a model pre-trained primarily on English corpus as English-centric model.

081 an English-centric model which has potentially  
082 superior English reasoning abilities and hope these  
083 can be transferred to other languages? Or should  
084 we start with the multilingual models which are  
085 generally assumed to have better multilingual  
086 transferability, but may lag behind in English  
087 reasoning skills?

088 In this paper, we investigate the following three  
089 research questions:

- 090 • How does the backbone (e.g., a multilin-  
091 gual pre-trained model or an English-centric  
092 model) affect multilingual reasoning?
- 093 • How does the source language used for down-  
094 stream task finetuning affect multilingual  
095 reasoning on other target languages? For  
096 example, will English always be the most  
097 effective source language for English-centric  
098 models?
- 099 • How does task type affect multilingual rea-  
100 soning, e.g., will the reasoning ability be  
101 transferred better across languages in some  
102 reasoning tasks?

103 To answer these questions, we consider four  
104 tasks that require distinct types of reasoning,  
105 namely Natural Language Inference, Logical  
106 Reasoning, Commonsense Reasoning, and Arith-  
107 metic Reasoning, and three popular multilingual  
108 and English-centric models, i.e., BLOOM (Scao  
109 et al., 2022), Pythia (Biderman et al., 2023) and  
110 LLaMA (Touvron et al., 2023a). We conduct exten-  
111 sive experiments in these multilingual downstream  
112 tasks, and have the following key observations:

- 113 • The multilingual pre-trained model does not  
114 always outperform an English-centric model,  
115 especially for languages seen or rarely seen  
116 for both models. For instance, LLaMA  
117 achieves a maximum of 9.9% and a minimum  
118 of 0.54% more average accuracy gain than  
119 BLOOM on Turkish and Greek, respectively,  
120 both are rarely seen for the two models (§3.2);
- 121 • Incorporating a small amount of multilingual  
122 data during the pre-training stage can have a  
123 significant impact on English-centric models.  
124 For example, though LLaMA is trained on  
125 French and Spanish data with a size of  
126 approximately 50 times less than BLOOM,  
127 it still outperforms BLOOM by up to 23% on  
128 these languages (§3.2);

- The choice of language utilized during fine-  
tuning becomes less important when the  
English-centric model scales up (§3.3);
- Different types of tasks show different multi-  
lingual transfer abilities, e.g., logical reason-  
ing knowledge can be transferred better across  
languages than others. However, as the model  
size increases, this gap tends to narrow (§3.4).

The experiment code is publicly available  
to promote reproducibility and facilitate further  
research.<sup>2</sup>

## 2 Language Versatilists and Specialists

In this section, we describe multilingual pre-  
training, with a focus on the curse of multilin-  
gual pretraining, and then discuss English-centric  
pretraining, with a series of evidence to show  
the potential of English-centric model possessing  
multilingual transfer ability.

**Multilingual pre-training** Multilingual pre-  
training offers a straightforward way to create  
language versatilists (Devlin et al., 2019; Conneau  
et al., 2018a; Xue et al., 2021; Shliazhko et al.,  
2022; Lin et al., 2021; Scao et al., 2022). The main  
idea is to combine monolingual corpora in different  
languages, upsampling those with less data, and  
training a regular language model on the combined  
data. After learning multiple languages that use  
diverse scripts and belong to various language  
families, the models are expected to possess  
*cross-lingual transfer ability*, i.e., the model can  
generalize to target languages (Pires et al., 2019;  
Wu and Dredze, 2019; Hu et al., 2020; Zhu et al.,  
2023) when downstream labeled training data  
is only available in the source language, which  
is especially important for low-resource target  
languages (Conneau et al., 2018a).

**Curse of multilingual pre-training** Conneau  
et al. (2018a) demonstrated that including more  
languages in a single model can improve per-  
formance for low-resource languages but hurt  
performance for high-resource languages. Fur-  
thermore, Wang et al. (2020) shows that negative  
interference between languages also leads to  
degraded performance on low-resource languages.  
As such, prior work had to find a trade-off between  
supporting more languages and obtaining better  
performance on a certain set of languages, such as

<sup>2</sup>URL is anonymized pending the reviewing process.

Language	Script	BLOOM	LLaMA
English (EN)	Latin	0.485	~4.666
Chinese (ZH)	ZH-ideograms	0.261	-
French (FR)	Latin	0.208	~0.004
Spanish (ES)	Latin	0.175	~0.004
Arabic (AR)	Arabic	0.075	-
Vietnamese (VI)	Latin	0.043	-
Hindi (HI)	Devanagari	0.025	-
Urdu (UR)	Perso-Arabic	0.003	-
Swahili (SW)	Latin	<0.001	-
Bulgarian (BG)	Cyrillic	-	~0.004
Russian (RU)	Cyrillic	-	~0.004
German (DE)	Latin	-	~0.004
Turkish (TR)	Latin	-	-
Greek (EL)	Greek	-	-
Thai (TH)	Brahmic	-	-

Table 1: Disk size (TB) of the pre-training data per language. 15 languages in the XNLI dataset are shown and sorted by their size in BLOOM. The numbers for LLaMA are roughly estimated based on Touvron et al. (2023a).

increasing model and vocabulary size (Conneau et al., 2018a; Wang et al., 2020), and learning additional language-specific parameters through adapters (Pfeiffer et al., 2022).

**English-centric pre-training** While only 13% of the world’s population speaks English, the vast majority of NLP research is done on English. Consequently, numerous models are pre-trained using a corpus that is primarily in English, while without explicitly excluding other languages during data collection (Brown et al., 2020; Chowdhery et al., 2022; Black et al., 2021; Wang and Komatsuzaki, 2021; Black et al., 2022; Biderman et al., 2023; Zhang et al., 2022; Touvron et al., 2023a). For example, English accounts for approximately 97.4% in the Pile (Gao et al., 2020), an 825GB dataset used by many pre-trained models (Black et al., 2021; Wang and Komatsuzaki, 2021; Black et al., 2022; Biderman et al., 2023), 93% in training data of GPT-3 (Brown et al., 2020), and around 99% in training data of LLaMA (Touvron et al., 2023a). In comparison, the largest constitution, i.e., English, only accounts for 30% in the ROOTS (Laurençon et al., 2022), which is the multilingual corpus for pretraining BLOOM (Scao et al., 2022). Table 1 compares the data size of the pretraining corpus for BLOOM and LLaMA model across 15 languages from XNLI dataset (Conneau et al., 2018b).

**Harbingers of multilingual transfer ability in English-centric models** Multiple lines of

evidence suggest that English-centric models have the potential for multilingual transfer capability. On the one hand, large English-centric models perform comparably with multilingual models on multilingual question-answering tasks (Chowdhery et al., 2022) and translating other languages into English (Brown et al., 2020; Chowdhery et al., 2022), though still lagging behind in translating into other languages. On the other hand, prior work suggests that the source of multilingual transfer ability may not be solely attributed to the multilingual pretraining process, as monolingual models also learn some abstractions that generalize across languages (Artetxe et al., 2020). High-level knowledge-transferring phenomena have been observed in other modalities, such as from English to Python (Hernandez et al., 2021), from ‘non-linguistic data with grammatical structure’ to language (Papadimitriou and Jurafsky, 2020; Ri and Tsuruoka, 2022), and from language to vision (Lu et al., 2021). Similarly, the presence of innate biological properties of the brain that constrain possible human languages was posited to explain why children learn languages so quickly despite the poverty of the stimulus (Chomsky, 1981; Legate and Yang, 2002).

## 3 Experiments

### 3.1 Setup

**Models** We consider both multilingual models and English-centric models and choose the three most popular models as the backbone in our experiments. The details of them are listed as follows:

- **BLOOM** (Scao et al., 2022): a series of models trained on ROOTS (Laurençon et al., 2022), a multilingual corpus containing 341 billion tokens from 46 natural languages and 13 programming languages. We consider three model sizes, i.e., 560M, 1.7B, and 7.1B, in our experiments;
- **Pythia** (Biderman et al., 2023): a family of models trained on the Pile (Gao et al., 2020), an English-centric corpus contains 207 billion tokens after deduplication. The overall number of tokens of the deduplicated Pile is on par with ROOTS. We consider three model sizes, i.e., 410M, 1.4B, and 6.5B, in our experiments;
- **LLaMA** (Touvron et al., 2023a): a series of models trained on various English-centric

Model	XNLI																
	en	zh	fr	es	ar	vi	hi	ur	sw	bg	ru	de	tr	el	th	Ave(15)	Ave(3)
<b>BLOOM-7.1B</b>	81.38	<b>70.72</b>	75.25	<b>77.96</b>	<b>69.46</b>	<b>69.96</b>	<b>62.75</b>	<b>57.33</b>	<b>56.25</b>	50.52	59.60	59.72	44.15	51.22	46.59	<b>62.19</b>	<b>75.78</b>
<b>Pythia-6.9B</b>	83.77	61.84	70.10	70.84	56.03	55.91	47.31	46.31	45.59	61.88	61.10	65.89	<b>54.39</b>	<b>61.50</b>	<b>51.42</b>	59.59	71.90
<b>LLaMA-6.7B</b>	<b>86.85</b>	61.82	<b>76.99</b>	77.56	52.69	54.71	46.97	45.51	40.58	<b>72.79</b>	<b>73.09</b>	<b>75.81</b>	51.12	57.39	46.57	61.36	75.22

Model	GSM8K				LogiQA				XCOPA			
	en	zh	fr	Ave(3)	en	zh	fr	Ave(3)	en	zh	fr	Ave(3)
<b>BLOOM-7.1B</b>	11.60	<b>8.80</b>	14.00	11.47	25.81	23.35	23.96	24.37	54.00	51.40	48.80	51.40
<b>Pythia-6.9B</b>	12.80	6.00	10.00	9.60	32.10	27.96	30.26	30.11	50.80	50.00	53.40	51.40
<b>LLaMA-6.7B</b>	<b>27.20</b>	7.20	<b>18.00</b>	<b>17.47</b>	<b>37.63</b>	<b>31.34</b>	<b>33.79</b>	<b>34.25</b>	<b>85.60</b>	<b>59.80</b>	<b>71.40</b>	<b>72.27</b>

Table 2: Accuracy of similar-sized multilingual and English-centric models on each test language after finetuning on English task data. The language is sorted by the pre-train data size in BLOOM as shown in Table 1. Ave(15) refers to the average results of all 15 test languages and Ave(3) is the average of the top three resourced languages (EN, ZH, FR) in BLOOM. Best result is in bold for each language. Full results of all model sizes and all the training languages are shown in the Appendix B.

corpus, summing up to tokens (1.4 trillion), much larger than that in ROOTS (341 billion) and the Pile (207 billion). Currently, LLaMAs are one of the most well-performed open-sourced models among similar-sized models. We consider three model sizes, i.e., 6.7B, 13B, and 32.5B, in our experiments.

**Datasets** We focus on multilingual reasoning ability in different models and consider four datasets that require distinct reasoning abilities, i.e., XNLI (Conneau et al., 2018b), LogiQA (Liu et al., 2021), XCOPA (Ponti et al., 2020) and GSM8K dataset (Cobbe et al., 2021). We create multilingual versions of a dataset through Google Translate API<sup>3</sup> if it doesn't have. We elaborate more details in Appendix A for each dataset.

**Implementation Details** We separately fine-tune the above 9 models on each language from the four datasets. As full fine-tuning becomes less feasible when the model gets larger, we adopt Low-Rank Adaptation (LoRA; (Hu et al., 2022)) and Int8 quantization (Dettmers et al., 2022) to perform compute and memory-efficient fine-tuning. With the above techniques, the finetuning and inference for the considered largest 32.5B model can be accomplished on a single NVIDIA A100-80GB GPU. Additionally, instead of using all the 400k training instances for each language in the XNLI dataset, we limit the number of training instances to 9k, with 3k for each class, to reduce computation. We set the batch size to 32, the learning rate to 3e-4, and the number of epochs to 3. We adopt instruction fine-tuning (Wei et al.,

2021; Sanh et al., 2021) instead of classifier-based fine-tuning (Devlin et al., 2019) for classification tasks, which injects certain abilities without adding additional modules. The number of instances and instruction templates for each dataset are listed in the Appendix Table 3. During inference, we compare the perplexity of each option to decide the label for classification tasks following (Brown et al., 2020), and we adopt the open-sourced OpenICL toolkit (Wu et al., 2023) for implementation. We always use English prompts as suggested by prior works (Lin et al., 2021; Muennighoff et al., 2022).

### 3.2 Findings for RQ1

*"RQ1: How does the backbone (e.g., a multilingual pre-trained model or an English-centric model) affect multilingual reasoning?"*

To facilitate the discussion, we use three models of similar parameters, i.e., BLOOM-7.1B, Pythia-6.9B, and LLaMA-6.7B. We begin by showing the overall accuracy of the three models on all the languages after training on English task data, as shown in Table 2. Then, we split the train languages into four categories based on the pretraining languages of BLOOM and LLaMA as listed in Table 1: (1) seen for both, (2) rarely seen for both, (3) seen for BLOOM but rarely for LLaMA, and (4) seen for LLaMA but rarely for BLOOM. We visualize the results in Figure 1, where the zero-shot accuracy is subtracted to better reflect performance gain brought from additional training on the certain source language.

**A minimal amount of multilingual data makes a lot in English-centric models** As shown in Table 2, LLaMA achieves comparable or better

<sup>3</sup><https://cloud.google.com/translate>

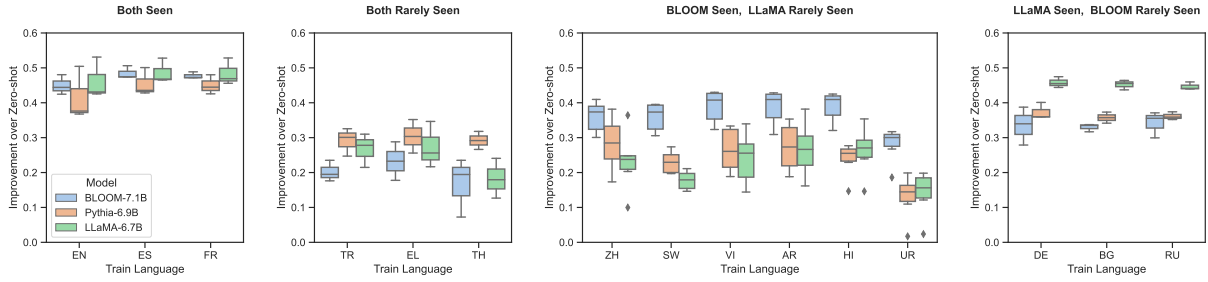


Figure 1: Evaluating BLOOM-7.1B and LLaMA-6.7B on four groups of languages, i.e., both seen during pre-training, both rarely seen during pre-training, seen for BLOOM but rarely seen for LLaMA, and seen for LLaMA but rarely seen for BLOOM. The zero-shot accuracy is subtracted to better reflect performance gain brought from additional training on the certain source language.

overall performance on the multilingual test sets, with an average accuracy of 61.39% compared to 62.19% of BLOOM. Even on languages frequently seen by BLOOM (i.e., EN, ZH, and FR), the average performance of LLaMA can still match (XNLI) or outperform (GSM8K, LogiQA, and XCOPIA) BLOOM. During pre-training, LLaMA only sees French and Spanish data with individual sizes equal to roughly 4 GB. By contrast, BLOOM has seen about 50 times data in these languages in the pre-training stage. Nevertheless, when evaluating LLaMA on French, the accuracy exceeds that of BLOOM by more than 1.7%, 4%, 10%, and 23% on XNLI, GSM8K, LogiQA and XCOPIA, respectively. LLaMA also achieves a very similar accuracy on Spanish, with BLOOM performing slightly better by a margin of 0.4% on XNLI.

However, for languages without any pre-training data (e.g., Chinese, Arabic, Vietnamese, etc.), the performance lags behind BLOOM by around 15% on XNLI but is still comparable or better on other tasks. Our findings suggest that incorporating a minimal amount of diverse low-resource language data during pre-training can result in a more capable multilingual pre-trained model, which outperforms models not trained on any data in those languages.

**LLaMA possesses better transfer ability across seen languages than BLOOM.** The first subplot of Figure 1 shows the accuracy improvements from directly zero-shot testing on the three languages seen by all models (i.e., EN, FR, ES) to first training on the three languages and then testing on these languages. LLaMA demonstrates better or comparable multilingual transfer ability for all the training languages. Since LLaMA was trained on mostly English texts, it is natural to

expect that it learns English data in finetuning better than multi-lingual models like BLOOM. This is consistent with the experimental result, where both minimum and maximum improvements for LLaMA are greater than those for BLOOM. Among the three models, Pythia has consistently lower improvements over zero-shot learning. We conjecture that the size of the English pre-training corpus has a positive correlation with a model’s multilingual transfer ability.

**Both English-centric models transfer better on rarely seen languages than BLOOM.** As illustrated in the second subplot of Figure 1, on the one hand, LLaMA exhibits more effective knowledge acquisition from Turkish (TR) and Greek (EL) data than BLOOM, which enhances its reasoning ability regardless of the language in which it is evaluated. This implies that a deep understanding of a single language could potentially enhance a model’s ability to comprehend unfamiliar languages more than a shallow understanding of multiple languages. On the other hand, Pythia emerges as the best-performing model when trained and evaluated on rarely seen languages by LLaMA and BLOOM. Considering the performance difference between Pythia and LLaMA, which are both English-centric models, we argue that the former’s superiority can partially be attributed to the different language distributions of their pre-training dataset excluding English data. This suggests that even with fewer overall pre-training data, models can have a better transfer result after pre-training in the specific language.

**Language coverage in pre-training is still important for multilingual transfer.** As illustrated in the third subplot of Figure 1, we found that BLOOM overall performed the best, surpassing

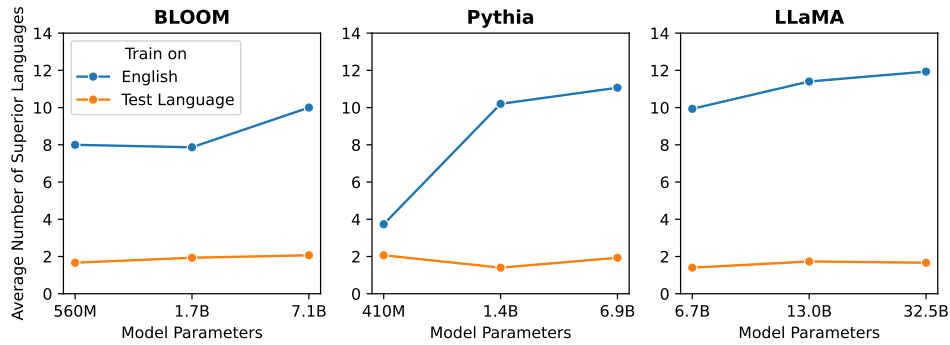


Figure 2: Average number of superior training languages compared with English and the test language.

the other two models by a great margin. This is not surprising because BLOOM is trained on them while others are rarely. Pythia comes as the second, with LLaMA being the last. The superiority of Pythia over LLaMA can be attributed to the difference in their pre-train datasets. For Pythia, its dataset consists of 97.4% of English data with the remaining for other languages, whereas for LLaMA, more than 99% of its pre-train data is in English. Therefore, we suspect that a slightly more diverse pre-train dataset in languages benefits Pythia towards capturing linguistic universals.

Finally, as illustrated in the fourth subplot of Figure 1, we show that when training and evaluating in languages that LLaMA has seen but BLOOM hasn't, the test accuracies of LLaMA are significantly higher than the other two models, with Pythia being the second. This further suggests language coverage in pre-training is important for both multilingual models and English-centric models.

### 3.3 Findings for RQ2

*"RQ2: How does the source language used for downstream task finetuning affect multilingual reasoning on other languages?"*

**For both multilingual and English-centric models, English appears to be a less suitable source language when the model scales up.** To investigate how the source language used for finetuning behaves on different models with different model sizes. We calculate the average number of superior source languages compared with English and the target language on the XNLI dataset. The value varies from 0 to 14, indicating the certain source language (i.e., English or the target language) is from the best to the worst among the total 15 languages, respectively. We show the

results in Figure 2.

As the model scales up, our experiments reveal that for all three models, there is a general increasing trend for the number of superior languages compared to English as model parameters grow. This observation can be attributed to the increasing capacity of the model, which enables it to capture more nuanced linguistic features. A possible explanation is as the increase of model capacity, the learning of other source languages becomes easier and consequently enhances the chances of identifying a more suitable source language other than English. These findings are applicable not only to multilingual models but also to both English-centric models.

**Training on target language may not be the best choice but can be a safe option.** While training in the target language is not always the optimal choice, we find it consistently yields good performance. Based on Figure 2, there are a small fraction of cases, with a number of approximately 2, where the accuracy difference is obtained by subtracting the accuracy of the model trained on each target language itself from trained on other languages, is positive. This finding suggests that incorporating target language data during training allows the model to better adapt to the specific characteristics of that language.

We further delve into each language to see if the on-average two superior languages are always the same for different models. To achieve this, we set the performance of the model trained on the target language as the baseline (0), and compute the relative performance gap of the model trained on each other source language. As shown in Figure 3, we find that such occurrences are primarily observed in Chinese (ZH), French (FR), Spanish (ES), and Urdu (UR), for LLaMA.

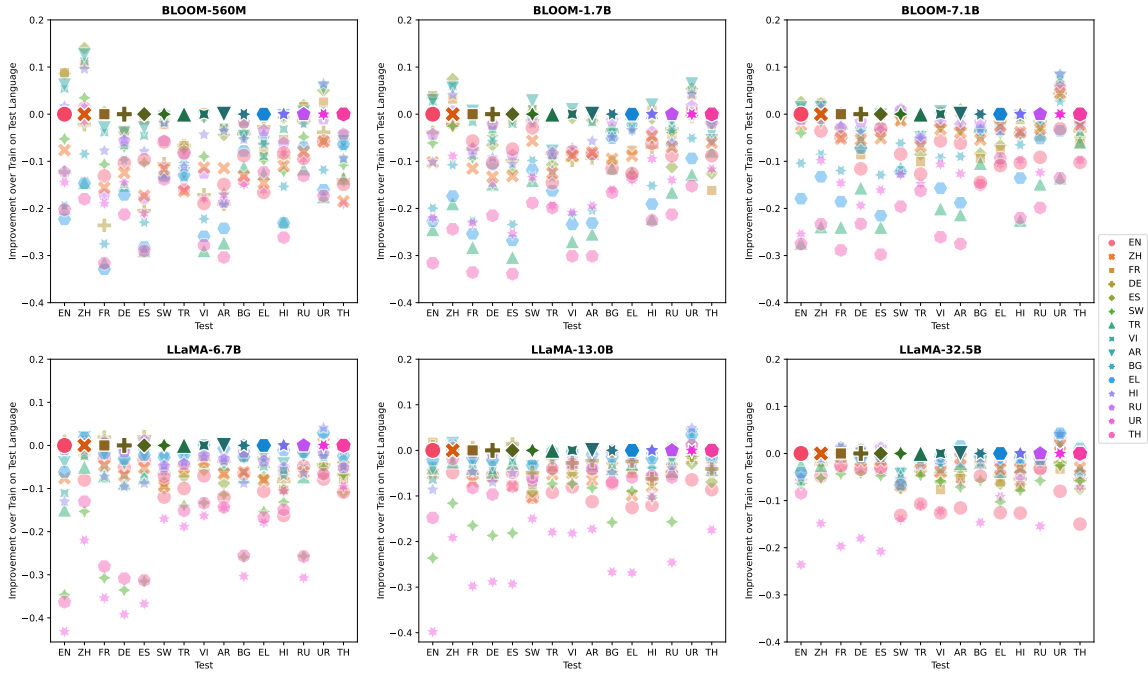


Figure 3: Accuracy gain of BLOOMs and LLaMAs on test languages by subtracting the performance of models trained on each test language from those trained on other languages.

473 While for BLOOM, they are mostly English (EN),  
 474 Chinese (ZH), and Urdu (UR). The results appear  
 475 to be complex as they are not highly correlated  
 476 with the language frequency observed during the  
 477 pre-training stage. For instance, LLaMA has seen  
 478 English, French, and Spanish, while BLOOM  
 479 has seen English and Chinese. One possible  
 480 explanation for this can be the distinctive language  
 481 scripts used in Chinese (Chinese ideograms) and  
 482 Urdu (Perso-Arabic), which may not be well-suited  
 483 for acquiring knowledge related to reasoning.

484 **Languages used in finetuning become increasingly irrelevant as an English-centric model scales up.** In terms of Figure 3, as the parameters of LLaMA grow, the distribution of Y-coordinates (i.e., accuracy improvements) becomes more concentrated around the line  $y = 0$ , which corresponds to training and testing on the same language. Through a comparison between the LLaMA-6.7B and LLaMA-32.5B models, we find that the larger model not only exhibited fewer negative outliers, which were mostly associated with SW, UR, and TH as train languages, but also demonstrates significant accuracy improvements for other languages. As a result, the difference in accuracy between training on the target language and training on other languages is reduced when the model gets larger. In contrast, we do not observe

501 a clear trend with BLOOM as the model size  
 502 increased from 560M to 7.1B. Additionally, we find  
 503 the results on Pythia as shown in Appendix Figure 7  
 504 to be less conclusive than those on LLaMA, and  
 505 we attribute this to both the model size and the  
 506 English-centric pre-training.

### 3.4 Findings for RQ3

507 *"RQ3: How does task type affect multilingual reasoning, e.g., will the reasoning ability be transferred better across languages in some reasoning tasks?"*

508 Previous work finds the transfer performance on  
 509 'lower-level' tasks (e.g., POS-tagging, dependency  
 510 parsing, and NER) to be better correlated with  
 511 the syntactic similarity between languages, while  
 512 'high-level' tasks (e.g., NLI and QA) rely more on  
 513 other factors such as the size of pretraining corpora  
 514 of the target language (Lauscher et al., 2020). We  
 515 are interested to see whether transfer performance  
 516 also differs in different high-level reasoning tasks.  
 517  
 518  
 519  
 520

521 **Logical reasoning knowledge can be transferred better across languages than others, and such transferability on most tasks can be enhanced by scaling model size, even with a fixed English-centric pretraining corpus.** To measure the multilingual reasoning transfer ability for different tasks, we calculate the performance gap between  
 522  
 523  
 524  
 525  
 526  
 527

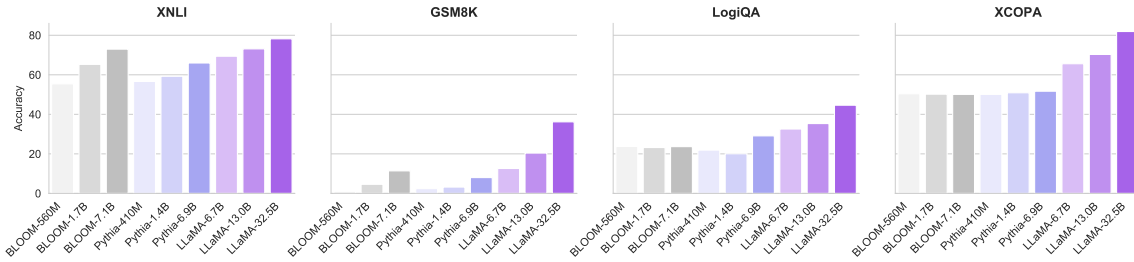


Figure 4: Average accuracy on other languages (i.e., FR and ZH) of each model trained on English task data across the four tasks.

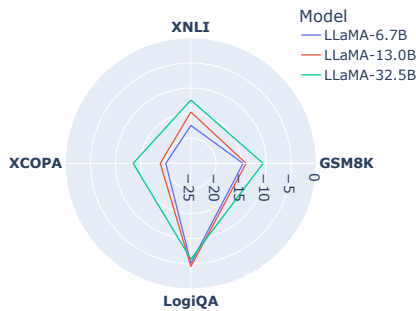


Figure 5: Performance gap between the average accuracy on other languages (i.e., FR and ZH) and English using the English-trained model. 0 refers to no performance gap, meaning the task ability transfers well from English to others.

the average accuracy on other languages and English using the English-trained model. We consider three languages, i.e., English, French, and Chinese for all the tasks. A value of 0 refers to no performance gap, meaning the reasoning ability transfers well from English to others. We show the results on LLaMA with various model sizes in Figure 5. The results indicate that LogiQA, which focuses on logical reasoning, exhibits the highest transferability across all the model sizes considered. On the other hand, XNLI, which tests with natural language inference, and GSM8K, which tests arithmetic reasoning, demonstrate comparatively lower levels of effectiveness. Furthermore, the figure indicates that increasing the model size generally leads to improved performance across most of the tasks, suggesting that multilingual reasoning transferability can be enhanced by increasing the model size, even if the training corpus remains constant. However, the results are fairly stable when the model scales up for LogiQA, with around 5% lower than the performance of testing on English, suggesting that solely increasing the model size only improves the

transfer ability to a certain amount.

### Multilingual pre-trained models fail on some multilingual reasoning tasks that English-centric models can handle.

We further study the multilingual reasoning transfer ability of different types of models on the four tasks. We show the average accuracy of the English-trained model when testing other languages in Figure 4. Notably, BLOOM-7.1B failed on the LogiQA dataset, exhibiting a level of performance that was no better than random guessing, while both Pythia-6.9B and LLaMA-6.7B, two English-centric models, achieves better performance. This suggests that a multilingual model may not possess sufficient capability to learn certain types of reasoning tasks as an English-centric model does. Additionally, both BLOOM-7.1B and Pythia-6.9B failed on the XCOPA dataset. In contrast, LLaMA-7B performed significantly better on both of these tasks, highlighting the importance of considering the fundamental capabilities of a language model in the context of multilingual reasoning tasks.

## 4 Conclusion

In this work, we investigate the multilingual transfer capabilities of both multilingual pre-trained and English-centric models, on four multilingual reasoning tasks. Our findings suggest that English-centric models possess significant multilingual transferability. We also found that English may not be the most effective source language for English-centric models, and different types of reasoning tasks exhibit varying multilingual transfer abilities. These findings offer practical insights for both pre-training and fine-tuning of the multilingual and English-centric models. We hope that our study will inspire further investigations and advancements in the development of more effective multilingual models.



## 590 Limitation

591 In this section, we discuss some potential lim-  
592 itations in our work. BLOOM and LLaMA,  
593 taken as representatives for language versatilis-  
594 t and specialist respectively, might not be strictly  
595 comparable because they were trained on different  
596 quantities of data. Hence, the results derived in  
597 our paper could tend to favor LLaMA which was  
598 pre-trained on more data considering all languages.  
599 To alleviate this inequality, we have conducted  
600 experiments on Pythia with a smaller pre-train  
601 dataset. If the corresponding result is still better  
602 than that of BLOOM, then we can conclude with  
603 stronger confidence that the specialist approach  
604 is superior. Nevertheless, noting that the quality  
605 of pre-train datasets can also vary, which makes  
606 Pythia and BLOOM still not strictly comparable.  
607 We acknowledge such possible deviations in the  
608 amount and quality of the pre-training corpus for  
609 the three models, and we recommend that future  
610 research pays more attention to it. In addition, we  
611 only evaluated the performance of supervised task  
612 fine-tuning in our study. In future work, it would be  
613 worthwhile to consider other learning paradigms  
614 such as in-context learning (Brown et al., 2020).

## 615 References

616 Mikel Artetxe, Sebastian Ruder, and Dani Yogatama.  
617 2020. On the cross-lingual transferability of monolingual  
618 representations. In *Proceedings of the 58th Annual*  
619 *Meeting of the Association for Computational Linguistics*,  
620 pages 4623–4637, Online. Association for Computational  
621 Linguistics.

622 Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie  
623 Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah  
624 Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward  
625 Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der  
626 Wal. 2023. Pythia: A suite for analyzing large language  
627 models across training and scaling.

628 Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony,  
629 Leo Gao, Laurence Golding, Horace He, Connor Leahy,  
630 Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai  
631 Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan  
632 Tow, Ben Wang, and Samuel Weinbach. 2022. GPT-NeoX-  
633 20B: An open-source autoregressive language model. In  
634 *Proceedings of the ACL Workshop on Challenges &*  
635 *Perspectives in Creating Large Language Models*.

636 Sid Black, Leo Gao, Phil Wang, Connor Leahy, and Stella  
637 Biderman. 2021. GPT-Neo: Large Scale Autoregressive  
638 Language Modeling with Mesh-Tensorflow. If you use this  
639 software, please cite it using these metadata.

640 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah,  
641 Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan,  
642 Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020.  
643 Language models are few-shot learners. *Advances in*  
644 *neural information processing systems*, 33:1877–1901.

Noam Chomsky. 1981. A naturalistic approach to language  
and cognition. *Cognition and Brain Theory*, 4(1):3–22. 645 646

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,  
Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul  
Barham, Hyung Won Chung, Charles Sutton, Sebastian  
Gehrmann, et al. 2022. Palm: Scaling language modeling  
with pathways. *arXiv preprint arXiv:2204.02311*. 647 648 649 650 651

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark  
Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert,  
Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021.  
Training verifiers to solve math word problems. *arXiv*  
*preprint arXiv:2110.14168*. 652 653 654 655 656

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*, pages 8440–8451, Online. Association for Computational Linguistics. 657 658 659 660 661 662 663 664

Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina  
Williams, Samuel Bowman, Holger Schwenk, and Veselin  
Stoyanov. 2018a. [XNLI: Evaluating cross-lingual sentence representations](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2475–2485, Brussels, Belgium. Association for Computational Linguistics. 665 666 667 668 669 670 671

Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina  
Williams, Samuel Bowman, Holger Schwenk, and Veselin  
Stoyanov. 2018b. [Xnli: Evaluating cross-lingual sentence representations](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2475–2485. 672 673 674 675 676 677

Tim Dettmers, Mike Lewis, Younes Belkada, and Luke  
Zettlemoyer. 2022. Llm.int8(): 8-bit matrix multiplication  
for transformers at scale. *arXiv preprint arXiv:2208.07339*. 678 679 680

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina  
Toutanova. 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)*, pages 4171–4186. 681 682 683 684 685 686 687

Leo Gao, Stella Biderman, Sid Black, Laurence Golding,  
Travis Hoppe, Charles Foster, Jason Phang, Horace He,  
Anish Thite, Noa Nabeshima, et al. 2020. The pile: An  
800gb dataset of diverse text for language modeling. *arXiv*  
*preprint arXiv:2101.00027*. 688 689 690 691 692

Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam  
McCandlish. 2021. Scaling laws for transfer. *arXiv*  
*preprint arXiv:2102.01293*. 693 694 695

Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi  
Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022.  
Lora: Low-rank adaptation of large language models. In  
*International Conference on Learning Representations*. 696 697 698 699

Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham  
Neubig, Orhan Firat, and Melvin Johnson. 2020. Xtreme:  
A massively multilingual multi-task benchmark for  
evaluating cross-lingual generalisation. In *International Conference on Machine Learning*, pages 4411–4421. PMLR. 700 701 702 703 704 705

706	Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, et al. 2022. The bigscience roots corpus: A 1.6 tb composite multilingual dataset. <i>Advances in Neural Information Processing Systems</i> , 35:31809–31826.	Ryokan Ri and Yoshimasa Tsuruoka. 2022. Pretraining with artificial language: Studying transferable knowledge in language models. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7302–7315.	767 768 769 770 771
712	Anne Lauscher, Vinit Ravishankar, Ivan Vulić, and Goran Glavaš. 2020. From zero to hero: On the limitations of zero-shot language transfer with multilingual Transformers. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4483–4499, Online. Association for Computational Linguistics.	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In <i>AAAI spring symposium: logical formalizations of commonsense reasoning</i> , pages 90–95.	772 773 774 775 776
719	Julie Anne Legate and Charles D Yang. 2002. Empirical re-assessment of stimulus poverty arguments. <i>The Linguistic Review</i> , 19(1-2):151–162.	Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, et al. 2021. Multitask prompted training enables zero-shot task generalization. In <i>International Conference on Learning Representations</i> .	777 778 779 780 781
722	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 Diab, Veselin Stoyanov, and Xian Li. 2021. Few-shot Learning with Multilingual Language Models.	Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. <i>arXiv preprint arXiv:2211.05100</i> .	782 783 784 785 786 787
729	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2021. Logiqa: a challenge dataset for machine reading comprehension with logical reasoning. In <i>Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence</i> , pages 3622–3628.	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2022. Language Models are Multilingual Chain-of-Thought Reasoners.	788 789 790 791 792
735	Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. 2021. Pretrained transformers as universal computation engines. <i>arXiv preprint arXiv:2103.05247</i> , 1.	Oleh Shliashko, Alena Fenogenova, Maria Tikhonova, Vladislav Mikhailov, Anastasia Kozlova, and Tatiana Shavrina. 2022. mGPT: Few-Shot Learners Go Multilingual.	793 794 795 796
738	Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M. Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2022. Crosslingual Generalization through Multitask Finetuning.	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	797 798 799 800 801
745	Isabel Papadimitriou and Dan Jurafsky. 2020. Pretraining on non-linguistic structure as a tool for analyzing learning bias in language models. <i>arXiv preprint arXiv:2004.14601</i> .	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	802 803 804 805 806
748	Jonas Pfeiffer, Naman Goyal, Xi Lin, Xian Li, James Cross, Sebastian Riedel, and Mikel Artetxe. 2022. Lifting the curse of multilinguality by pre-training modular transformers. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3479–3495, Seattle, United States. Association for Computational Linguistics.	Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <a href="https://github.com/kingoflolz/mesh-transformer-jax">https://github.com/kingoflolz/mesh-transformer-jax</a> .	807 808 809
756	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4996–5001, Florence, Italy. Association for Computational Linguistics.	Zirui Wang, Zachary C. Lipton, and Yulia Tsvetkov. 2020. On negative interference in multilingual models: Findings and a meta-learning treatment. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4438–4450, Online. Association for Computational Linguistics.	810 811 812 813 814 815
761	Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020. Xcopa: A multilingual dataset for causal commonsense reasoning. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2362–2376.	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. In <i>International Conference on Learning Representations</i> .	816 817 818 819 820
762		Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122.	821 822 823 824 825 826 827

828	Shijie Wu and Mark Dredze. 2019. <a href="#">Beto, bentz, becas: The surprising cross-lingual effectiveness of BERT</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 833–844, Hong Kong, China. Association for Computational Linguistics.	880
829		881
830		882
831		883
832		884
833		885
834		886
835	Zhenyu Wu, YaoXiang Wang, Jiacheng Ye, Jiangtao Feng, Jingjing Xu, Yu Qiao, and Zhiyong Wu. 2023. Openicl: An open-source framework for in-context learning. <i>arXiv preprint arXiv:2303.02913</i> .	887
836		
837		
838		
839	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. <a href="#">mT5: A massively multilingual pre-trained text-to-text transformer</a> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 483–498, Online. Association for Computational Linguistics.	888
840		889
841		890
842		891
843		892
844		893
845		894
846		895
847	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> .	896
848		
849		
850		
851		
852	Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian Huang. 2023. Multilingual machine translation with large language models: Empirical results and analysis. <i>arXiv preprint arXiv:2304.04675</i> .	897
853		898
854		
855		
856		
857	<b>A Datasets and Templates</b>	
858	We consider the following four types of tasks that require distinct reasoning abilities:	
859		
860	• <b>Natural Language Inference:</b> we use XNLI dataset ( <a href="#">Conneau et al., 2018b</a> ), which is created by crowd-translating the dev and test portions of the English Multi-NLI dataset ( <a href="#">Williams et al., 2018</a> ) into 14 languages (French (fr), Spanish (ES), German (DE), Greek (EL), Bulgarian (BG), Russian (RU), Turkish (TR), Arabic (AR), Vietnamese (VI), Thai (TH), Chinese (ZH), Hindi (HI), Swahili (SW), and Urdu (UR));	
861		
862		
863		
864		
865		
866		
867		
868		
869		
870	• <b>Logical Reasoning:</b> we adopt LogiQA dataset ( <a href="#">Liu et al., 2021</a> ), which is sourced from expert-written questions for testing human logical reasoning. As the training set is only available in English and Chinese, we further translate both training and test splits into French with Google Translate API <sup>4</sup> ;	
871		
872		
873		
874		
875		
876		
877	• <b>Commonsense Reasoning:</b> we choose XCOPA dataset ( <a href="#">Ponti et al., 2020</a> ), which is a causal commonsense reasoning task in	
878		
879		
	which a model is given a premise sentence and must determine either the cause or effect of the premise from two possible choices. Since the dataset only provides multilingual test sets, we utilize the training set from the original English COPA release ( <a href="#">Roemmele et al., 2011</a> ) and translate it into Chinese and French with Google Translate API;	
	• <b>Arithmetic Reasoning:</b> we use GSM8K dataset ( <a href="#">Cobbe et al., 2021</a> ), which contains linguistically diverse grade school math word problems. <a href="#">Shi et al. (2022)</a> construct a multilingual test set which we directly adopt for our test set. To construct a multilingual training set, we further translate the English training set into French and Chinese with Google Translate API.	
	We show the number of instances and the template used in each dataset in Table 3.	
	<b>B Detailed Results</b>	
	The detailed results for the three BLOOM, Pythia, and LLaMA models across 15 languages on the XNLI dataset are shown in Table 4, Table 5, and Table 6, respectively. The results on other three datasets (i.e., GSM8K, XCOPA and LogiQA) are listed in Table 7.	
	<b>C Additional Figures</b>	
	We show the accuracy gain of BLOOMs and LLaMAs on test languages by subtracting the performance of models trained on each test language from those trained on other languages in Figure 6. This figure is complementary to Figure 3 which only shows the results for BLOOMs and LLaMAs in the paper. Similarly, we also show the accuracy gain by subtracting the performance of models trained on English from those trained on other languages in Figure 7. This figure corresponds to the average number of superior training languages compared with English in Figure 2 of the paper, and shows specifically which languages are better used for training given a test language.	

<sup>4</sup><https://cloud.google.com/translate>

	<b>#Train/#Test</b>	<b>Template</b>
<b>XNLI</b>	9,000/5010	Question: {premise} Based on the previous passage, is it true that "{hypothesis}"? Yes, No, or Maybe?  Answer: {output}
<b>XCOPA</b>	400/500	Question: {premise} Based on the previous passage, choose the most reasonable {cause   effect}. A: {choice1} B: {choice2}  Answer: {output}
<b>LogiQA</b>	7,376/651	Question: {context} {question} A: {choice1} B: {choice2} C: {choice3} D: {choice4}  Answer: {output}
<b>GSM8K</b>	7,473/250	Question: {input}  Answer: {output}

Table 3: Number of training and test instances for each dataset, as well as the templates used during fine-tuning and inference.







Model	Train	XCOPA				LogiQA				GSM8K			
		en	fr	zh	Average	en	fr	zh	Average	en	fr	zh	Average
BLOOM-560M	Zero-shot	50.00	50.00	50.00	50.00	20.28	20.28	20.28	20.28	2.40	2.00	1.20	1.87
	Average	48.73	50.13	51.13	50.00	22.63	24.37	22.43	23.14	3.07	1.73	1.33	2.04
	en	49.00	52.00	48.80	49.93	23.50	25.19	22.27	23.66	4.80	2.00	0.00	2.27
	fr	48.80	51.40	51.40	50.53	23.20	25.35	20.43	22.99	2.00	1.60	0.80	1.47
	zh	48.40	47.00	53.20	49.53	21.20	22.58	24.58	22.79	2.40	1.60	3.20	2.40
BLOOM-1.7B	Zero-shot	49.20	49.80	50.00	49.67	19.97	20.58	20.28	20.28	1.60	2.40	2.40	2.13
	Average	49.13	50.20	51.00	50.11	25.14	26.11	22.32	24.53	4.40	4.27	4.27	4.31
	en	48.20	50.20	50.20	49.53	25.65	25.35	21.04	24.01	5.60	5.20	4.00	4.93
	fr	49.60	49.80	50.60	50.00	26.73	28.11	22.73	25.86	4.00	4.80	2.80	3.87
	zh	49.60	50.60	52.20	50.80	23.04	24.88	23.20	23.71	3.60	2.80	6.00	4.13
BLOOM-7.1B	Zero-shot	49.80	51.60	50.00	50.47	23.04	20.89	19.97	21.30	2.80	3.20	2.40	2.80
	Average	52.00	49.33	50.67	50.67	26.16	26.73	24.83	25.91	11.07	12.00	8.27	10.44
	en	54.00	48.80	51.40	51.40	25.81	23.96	23.35	24.37	11.60	14.00	8.80	11.47
	fr	50.20	49.00	50.40	49.87	26.27	28.73	25.81	26.93	11.20	10.00	7.60	9.60
	zh	51.80	50.20	50.20	50.73	26.42	27.50	25.35	26.42	10.40	12.00	8.40	10.27
Pythia-410M	Zero-shot	50.00	49.80	50.00	49.93	20.28	23.50	20.28	21.35	2.80	2.00	3.20	2.67
	Average	50.13	50.80	48.93	49.96	23.76	21.76	21.81	22.44	2.53	2.27	2.00	2.27
	en	50.20	50.40	49.80	50.13	25.65	22.12	21.66	23.14	2.40	2.80	2.00	2.40
	fr	50.20	50.80	50.20	50.40	25.19	21.66	21.66	22.84	2.80	2.80	1.20	2.27
	zh	50.00	51.20	46.80	49.33	20.43	21.51	22.12	21.35	2.40	1.20	2.80	2.13
Pythia-1.4B	Zero-shot	50.00	50.00	50.00	50.00	20.28	20.89	20.28	20.48	2.00	2.00	1.60	1.87
	Average	49.73	50.73	49.53	50.00	22.63	21.30	21.61	21.85	6.27	4.00	4.67	4.98
	en	49.60	51.60	50.20	50.47	21.20	20.12	19.97	20.43	8.40	3.20	3.20	4.93
	fr	49.80	51.00	49.80	50.20	25.04	19.51	23.50	22.68	6.80	8.00	1.60	5.47
	zh	49.80	49.60	48.60	49.33	21.66	24.27	21.35	22.43	3.60	0.80	9.20	4.53
Pythia-6.9B	Zero-shot	50.00	50.40	50.00	50.13	21.97	22.27	20.28	21.51	4.80	3.20	2.00	3.33
	Average	50.33	51.93	49.20	50.49	28.21	27.96	26.01	27.39	10.27	8.67	7.07	8.67
	en	50.80	53.40	50.00	51.40	32.10	30.26	27.96	30.11	12.80	10.00	6.00	9.60
	fr	50.00	51.00	49.20	50.07	25.65	27.19	24.42	25.76	10.80	11.60	4.40	8.93
	zh	50.20	51.40	48.40	50.00	26.88	26.42	25.65	26.32	7.20	4.40	10.80	7.47
LLaMA-6.7B	Zero-shot	54.40	51.00	52.00	52.47	21.97	24.88	22.27	23.04	4.00	3.20	3.20	3.47
	Average	72.00	63.67	54.07	63.24	33.59	32.10	29.03	31.58	22.93	18.00	10.13	17.02
	en	85.60	71.40	59.80	72.27	37.63	33.79	31.34	34.25	27.20	18.00	7.20	17.47
	fr	72.20	65.40	51.00	62.87	37.79	36.41	33.03	35.74	24.40	21.60	7.20	17.73
	zh	58.20	54.20	51.40	54.60	25.35	26.11	22.73	24.73	17.20	14.40	16.00	15.87
LLaMA-13.0B	Zero-shot	62.20	52.20	50.60	55.00	25.35	26.42	20.28	24.01	5.20	3.20	3.20	3.87
	Average	85.40	76.33	61.40	74.38	38.91	37.22	34.66	36.93	30.93	25.73	16.93	24.53
	en	89.20	76.80	63.80	76.60	39.78	37.63	33.03	36.82	34.40	27.60	13.20	25.07
	fr	83.00	75.20	58.80	72.33	40.40	36.87	35.18	37.48	31.60	30.40	15.20	25.73
	zh	84.00	77.00	61.60	74.20	36.56	37.17	35.79	36.51	26.80	19.20	22.40	22.80
LLaMA-32.5B	Zero-shot	50.00	50.00	50.00	50.00	20.58	27.19	21.20	22.99	15.20	10.00	3.20	9.47
	Average	95.13	91.27	76.07	87.49	49.51	46.80	43.57	46.63	46.80	43.60	29.87	40.09
	en	95.40	90.00	73.80	86.40	50.54	45.93	43.32	46.59	46.80	46.40	26.00	39.73
	fr	95.40	93.00	77.20	88.53	51.15	47.93	41.32	46.80	51.20	45.60	28.40	41.73
	zh	94.60	90.80	77.20	87.53	46.85	46.54	46.08	46.49	42.40	38.80	35.20	38.80

Table 7: Detailed results of BLOOM, Pythia, and LLaMA on XCOPA, LogiQA, and GSM8K datasets.



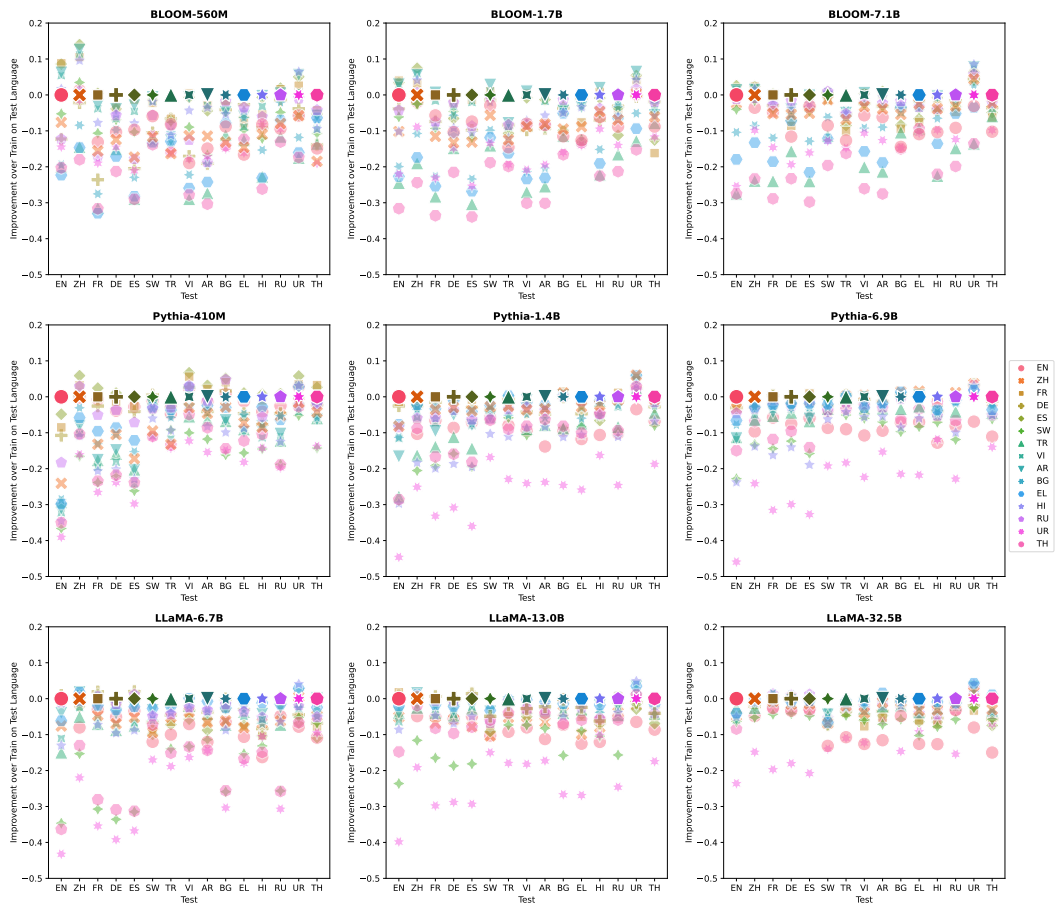


Figure 6: Accuracy gain of BLOOMs and LLaMAs on test languages by subtracting the performance of models trained on each test language from those trained on other languages.

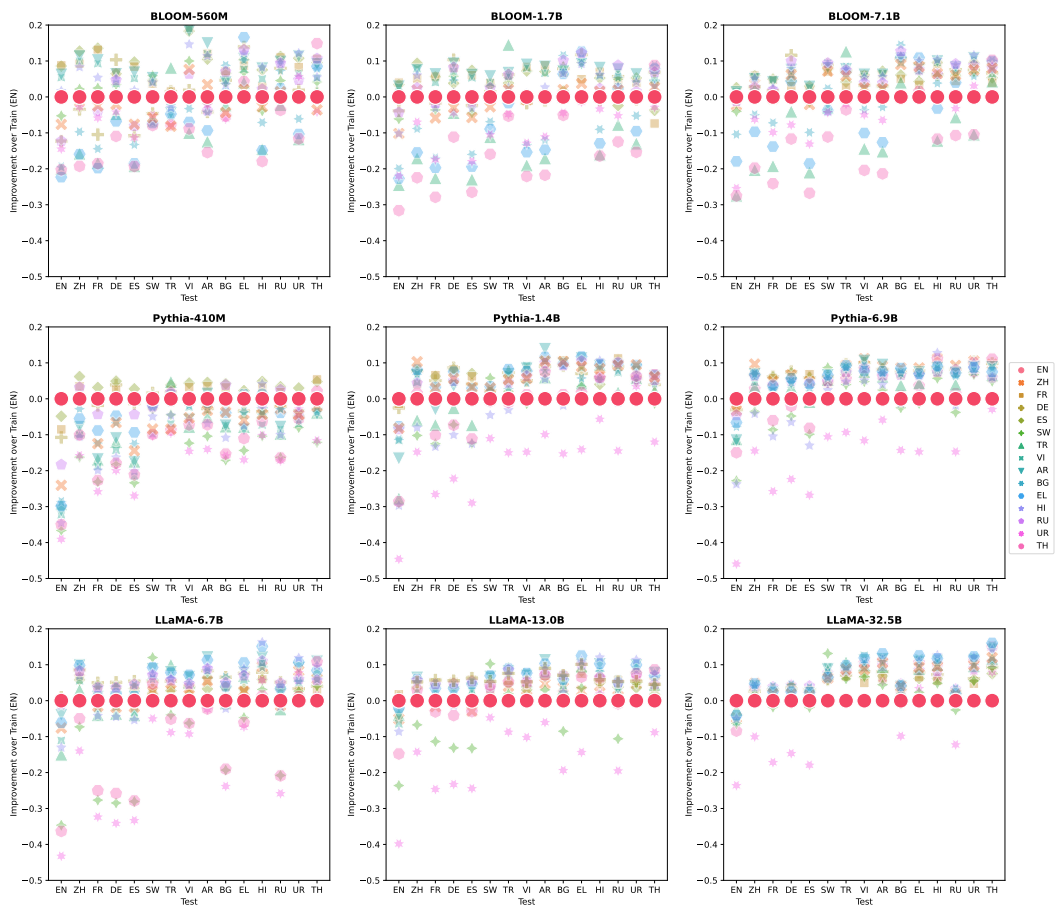


Figure 7: Accuracy gain of BLOOMs and LLaMAs on test languages by subtracting the performance of models trained on English from those trained on other languages.