

# A Recipe of Parallel Corpora Exploitation for Multilingual Large Language Models

Anonymous ACL submission

## Abstract

Recent studies have highlighted the potential of exploiting parallel corpora to enhance multilingual large language models, improving performance in both bilingual tasks, e.g., machine translation, and general-purpose tasks, e.g., text classification. Building upon these findings, our comprehensive study aims to identify the most effective strategies for leveraging parallel corpora. We investigate the impact of parallel corpora quality and quantity, training objectives, and model size on the performance of multilingual large language models enhanced with parallel corpora across diverse languages and tasks. Our analysis reveals several key insights: (i) filtering noisy translations is essential for effectively exploiting parallel corpora, while language identification and short sentence filtering have little effect; (ii) even a corpus containing just 10K parallel sentences can yield results comparable to those obtained from much larger datasets; (iii) employing only the machine translation objective yields the best results among various training objectives and their combinations; (iv) larger multilingual language models benefit more from parallel corpora than smaller models due to their stronger capacity for cross-task transfer. Our study offers valuable insights into the optimal utilization of parallel corpora to enhance multilingual large language models, extending the generalizability of previous findings from limited languages and tasks to a broader range of scenarios.

## 1 Introduction

Recent multilingual large language models (mLLMs), represented by BLOOM (Scao et al., 2022), MaLA500 (Lin et al., 2024b), and Aya (Üstün et al., 2024), have shown impressive capacity on diverse tasks across languages. Parallel corpora have emerged as crucial resources for enhancing mLLMs, both for specific tasks, e.g., machine translation (Xu et al., 2023; Alves et al.,

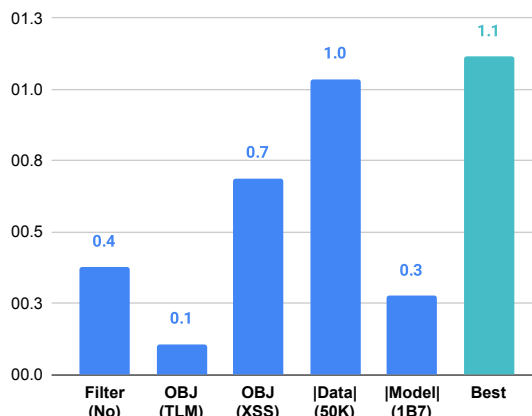


Figure 1: Average performance improvements achieved by mLLMs enhanced with parallel corpora compared to their base models. **Best: Continued pre-training of BLOOM-7B1 with the machine translation objective (MT) using 10K high-quality parallel sentences yields the best results.** Main variations explored include: **Filter (No)** (using the original data); **OBJ (TLM)** (translation language modeling objective); **OBJ (XSS)** (cross-lingual semantic similarity objective); **|Data| (50K)** (a larger 50K-sentence dataset); **|Modell| (1B7)** (BLOOM-1B7 model).

2024), and for general-purpose tasks (Cahyawijaya et al., 2023; Zhu et al., 2023; Li et al., 2023).

However, existing studies often overlook a comprehensive exploration of methodologies for harnessing parallel corpora. The quality and quantity of parallel corpora remain inadequately explored, inhibiting the full potential of such resources. Moreover, the influence of different training objectives and mLLM sizes across diverse languages and tasks remains under-investigated. This limitation impedes the generalization of parallel corpora exploitation methods across varied linguistic landscapes and task domains. Therefore, this paper aims to address these gaps by presenting a comprehensive recipe for exploiting parallel corpora for mLLMs. We focus on four key factors,

with some main results shown in Figure 1.

**Quality:** Kreutzer et al. (2022) highlight large portions of low-quality data within current massive parallel corpora. We explore three dimensions of quality: translation accuracy, sentence length, and language identification. Our findings emphasize the critical role of translation quality in exploiting parallel corpora, while showing that sentence length filtering and language identification have minimal impact.

**Quantity:** Acquiring substantial amounts of high-quality parallel corpora presents a significant challenge, especially for low-resource languages. Our study examines the minimum corpus size necessary to achieve performance improvements across diverse tasks. Remarkably, we find that even a corpus of just 10K sentences can yield results comparable to those obtained from much larger datasets.

**Objective:** Previous studies (Cahyawijaya et al., 2023) have investigated the effectiveness of different training objectives and their combinations on classification tasks of Indonesian local languages, using smaller-sized mLLMs up to 1B7. We extend this investigation by examining the impact of various training objectives and their combinations on larger mLLMs across a range of languages and tasks. Our experiments demonstrate that employing the machine translation objective produces the most promising results.

**Model Size:** The size of mLLMs can greatly impact their ability to comprehend instructions derived from parallel corpora. Our findings indicate that larger mLLMs exhibit superior comprehension and cross-task transferability compared to their smaller counterparts. Consequently, they achieve more substantial improvements across a broader spectrum of tasks.

In light of the critical role parallel corpora play in mLLMs, our study provides a comprehensive recipe for effectively exploiting parallel corpora. We have identified four primary factors: quality (§4), quantity (§5), objective (§6), and model size (§7). Our detailed analysis of these factors reveals their great impact on mLLM performance across diverse languages and tasks. By delving into these aspects, we offer actionable insights that can inform the development and optimization of strategies for parallel corpora exploitation, ultimately contributing to the advancement of mLLMs in both bilingual and general-purpose tasks.

## 2 Related Work

### 2.1 Parallel Data for Multilingual Language Models

Over the years, multilingual language models have evolved from earlier, smaller models, such as XLM (Conneau and Lample, 2019), XLM-R (Conneau et al., 2020), and Glot500 (Imani et al., 2023), to more recent, larger models, including BLOOM (Scao et al., 2022), MaLA500 (Lin et al., 2024b), and Aya (Üstün et al., 2024). These models consistently demonstrate strong performance across various downstream tasks (Ahuja et al., 2023; Lin et al., 2024c).

Parallel corpora have played a pivotal role in both the analysis (Piqueras and Søgaard, 2022; Lin et al., 2024a) and enhancement (Conneau and Lample, 2019; Ouyang et al., 2020; Yang et al., 2020; Huang et al., 2019; Chi et al., 2021a; Wei et al., 2021; Hu et al., 2021; Chi et al., 2021b; Reid and Artetxe, 2022b; Liu et al., 2023) of small multilingual language models.

In the era of mLLMs, parallel corpora are constructed as instruction data and used to enhance mLLMs through supervised fine-tuning (Cahyawijaya et al., 2023; Zhu et al., 2023; Li et al., 2023). Specifically, Cahyawijaya et al. (2023) propose three methods of incorporating parallel corpora as instruction tuning data: Machine Translation (MT), Translation Language Modeling (TLM), Cross-Lingual Semantic Similarity (XSS) (see §3.3). However, their evaluation is limited to small models with up to 1.7 billion parameters and focuses solely on classification tasks within Indonesian local languages. Both Zhu et al. (2023) and Li et al. (2023) propose using machine-translation-style instruction data to improve mLLMs but do not explore different training objectives. While these studies yield promising results, their scope is limited. Firstly, they fail to explore critical factors such as the quality and quantity of parallel corpora, considering the high cost of collecting high-quality and massive parallel corpora, especially for low-resource languages. Secondly, their investigations do not encompass an in-depth analysis of training objectives and mLLMs with varied model sizes across diverse languages and tasks.

### 2.2 Key Elements for Language Modeling

Previous research has extensively examined critical factors essential for the pretraining and enhancement of language models.

**Quality:** Kreutzer et al. (2022) conducted manual audits of prevalent monolingual and parallel corpora, revealing significant portions of low-quality data, particularly in corpora for low-resource languages. Follow-up studies have investigated the impact of data quality on model performance. Artetxe et al. (2022) observed that similar results on downstream tasks can be achieved regardless of the degree of quality of the corpus used for pretraining, while other studies found that the quality of parallel corpora matters for machine translation (Ranathunga et al., 2024) and general-purpose tasks (Reid and Artetxe, 2022a).

**Quantity:** Recent works (Chen et al., 2023; Zhou et al., 2023; Gupta et al., 2023) have focused on the impact of fine-tuning with small amounts of high-quality instruction data, such as one or a few thousand instances, showing promising performance gains in evaluation tasks. Xu et al. (2023) demonstrate that as few as 10K high-quality parallel sentences can significantly enhance machine translation performance.

**Objective:** Different training objectives based on parallel corpora for enhancing mLLMs can be viewed as distinct instructions. Wang et al. (2023) explore the impact of various types of instruction tuning data and find that their combination can be optimal in certain scenarios.

**Model Size:** Recent studies indicate that scaling up language models enhances their capability to excel in diverse and complex reasoning tasks (Wei et al., 2022, 2023; Lu et al., 2023). Follow-up studies (Shu et al., 2023; Wei et al., 2023) further illustrate distinct behavioral differences between larger and smaller models.

However, these factors have not yet been comprehensively explored in the context of leveraging parallel corpora to enhance mLLMs across diverse languages and tasks.

## 3 Setup

### 3.1 Language

We use two criteria for language selection. Firstly, we select languages well covered by mLLMs and evaluation benchmarks, allowing for robust evaluation across diverse downstream tasks. Secondly, we select typologically diverse languages, enabling our investigation to generalize to a wide range of low-resource languages. Thus, we select five languages: Arabic (ar), Spanish (es), Hindi (hi), Vietnamese (vi) and Chinese (zh).

### 3.2 Data

We utilize the OPUS100 dataset (Zhang et al., 2020), an English-centric multilingual corpus, to gather parallel sentences between English (en) and each target language. The quality of OPUS100 is assessed across three dimensions:

**Translation Quality** Manual quality assessment of the vast amount of parallel corpora is impractical. Instead, we employ COMETKIWI (Rei et al., 2022)<sup>1</sup>, a tool for estimating the quality of machine translation outputs across multiple languages. We set a COMETKIWI score threshold  $\tau_c$ , retaining parallel corpora with scores not lower than  $\tau_c$ .

**Sentence Length** Given the variation in character length across languages, we avoid using it as a metric for consistency. Instead, we measure sentence length by the number of tokens, as determined by the tokenizer of our chosen mLLM, BLOOM-7B1. We establish a length threshold  $\tau_l$ , retaining parallel corpora where both source and target sentences contain no fewer than  $\tau_l$  tokens.

**Language Identification** To identify sentences potentially not in the correct language, we employ GlotLID (Kargaran et al., 2023), an open-source language identification model. This language identification filter is applied to both the source and target sentences.

### 3.3 Training

We select the BLOOM series (Scao et al., 2022) for our investigation due to its offering of different sizes of mLLMs pretrained for the five target languages under consideration. We explore BLOOM models of various sizes, including 7B1, 3B, and 1B7. Due to limited computational resources, we continue pretraining BLOOM using LoRA (Hu et al., 2022), which is known for its competitive performance compared to full-parameter training. We configure the learning rate to  $1e-4$ , weight decay to 0.1, and set the rank of LoRA to 16. The maximum sequence length for both source and target sentences is set to 128. To maintain consistency across experiments with different quantities of parallel corpora, we ensure a uniform training budget of 50K parallel sentences. Specifically, we calculate the number of epochs as 50K divided by the number of sentences considered from the

<sup>1</sup><https://huggingface.co/Unbabel/wmt23-cometkiwi-da-xxl>

| Objective | Template  |
|-----------|---|
| MT        | Translate the following text from [SOURCE_LANG] to [TARGET_LANG].<br>Text: [SOURCE_TEXT]<br>Translation: [TARGET_TEXT]                |
| TLM       | [INPUT_TEXT]. Denoise the previous [TARGET_LANG] text to its equivalent sentence in [SOURCE_LANG]: [SOURCE_TEXT]<br>[TARGET_TEXT]     |
| XSS       | [SOURCE_LANG] sentence: [SOURCE_TEXT]<br>[TARGET_LANG] sentence: [TARGET_TEXT]<br>Do the two sentences have the same meaning? [LABEL] |

Table 1: Templates of MT, TLM, and XSS for instruction data construction based on parallel corpora.

OPUS100 dataset. The batch size is 128, and we save the checkpoints every 20 steps.

Following Cahyawijaya et al. (2023), we construct the data for instruction tuning based on the parallel corpora by three distinct patterns: Machine Translation (MT), Translation Language Modeling (TLM), and Cross-Lingual Semantic Similarity (XSS). Table 6 presents the templates for these three objectives. Here, [SOURCE\_LANG] and [TARGET\_LANG] represent the language names of the source and target languages, respectively. In our study, we consider both English-to-target-language and target-language-to-English directions, where [SOURCE\_LANG] represents English or [TARGET\_LANG] represents English. For MT, [SOURCE\_TEXT] and [TARGET\_TEXT] refer to the parallel sentences in the source and target languages, respectively. For TLM, a portion of tokens in [TARGET\_TEXT] are masked to generate [INPUT\_TEXT]. For XSS, our objective is to predict whether parallel sentences [SOURCE\_TEXT] and [TARGET\_TEXT] are semantically similar, with [LABEL] being “Yes” or “No”. Specifically, we utilize the parallel corpora as positive examples and introduce perturbations to [TARGET\_TEXT] to construct negative examples. We consider applying the objectives both individually and in combination.

### 3.4 Evaluation

We conduct evaluation across five diverse benchmarks: FLORES (Costa-jussà et al., 2022), MUSE (Lample et al., 2018), MLQA (Lewis et al., 2020), XQuAD (Artetxe et al., 2020), and SIB (Adelani et al., 2023). A comprehensive overview of these benchmarks is available in Table 2. Our evaluation spans both classification tasks (SIB) and generation tasks (FLORES, MUSE, MLQA, and XQuAD), covering a spectrum of cross-language (FLORES, MUSE, and MLQA) and in-language tasks (XQuAD and SIB).

| Dataset | Task                | IDatal      | Metric    | I/C | C/G |
|---------|---------------------|-------------|-----------|-----|-----|
| FLORES  | Machine Translation | 1012        | COMETKIWI | C   | G   |
| MUSE    | Word Translation    | 1500        | F1        | C   | G   |
| MLQA    | Question Answering  | 4918 - 5495 | F1        | C   | G   |
| XQuAD   | Question Answering  | 1190        | F1        | I   | G   |
| SIB     | Text Classification | 204         | Acc       | I   | C   |

Table 2: Details of evaluation benchmarks. IDatal: Number of samples for evaluation. I/C: In-language/Cross-language. C/G: Classification/Generation.

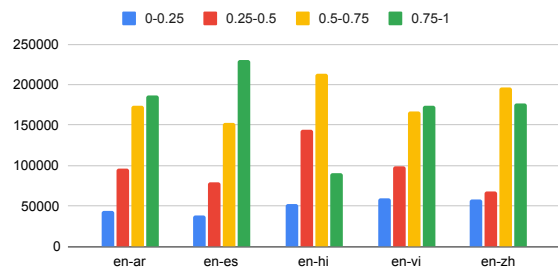


Figure 2: Translation quality measured by COMETWIKI of 500K parallel sentences from OPUS100 for our five language pairs. The COMETWIKI scores are segmented into four ranges: 0-0.25, 0.25-0.5, 0.5-0.75, and 0.75-1. Higher scores represent better translation quality.

For translation tasks within FLORES and MUSE, we explore bidirectional translation: from English to other languages (en-xx) and from other languages to English (xx-en). Additionally, for MLQA, we evaluate scenarios where questions are in English and the passages and answers are in other languages (en-xx), as well as situations where questions are in other languages and the passages and answers are in English (xx-en).

To provide a thorough understanding of our evaluation procedures, we offer detailed prompts for each task in §A. In all experiments, we employ a 2-shot in-context learning approach, where the model is given two examples appended to the input query to aid in making predictions.

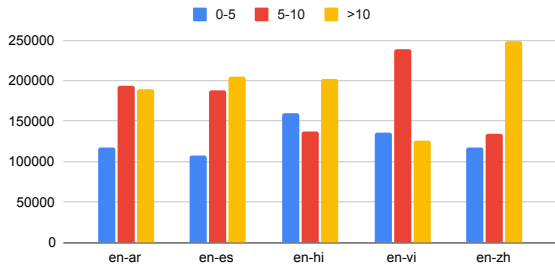


Figure 3: Sentence length of 500K parallel sentences from OPUS100 for our five language pairs. The three categories are 0-5, 5-10, greater than 10 tokens.

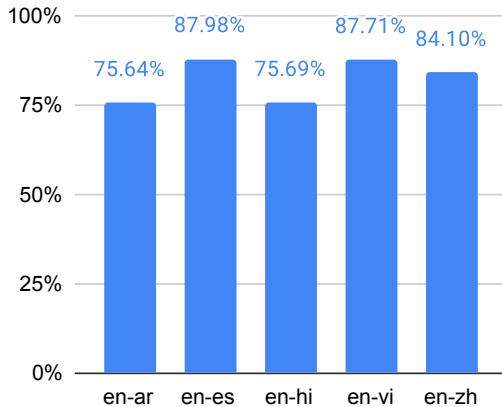


Figure 4: Percentage of sentences retained after language identification filtering of 500K parallel sentences from OPUS100 for our five language pairs.

## 4 Quality

### 4.1 Quality of OPUS100

We measure the quality of 500K parallel sentences from OPUS100 for our five language pairs using three key metrics: translation quality, sentence length, and language identification accuracy, as illustrated in Figures 2–4.

**A considerable portion of OPUS100 is of low-quality.** All quality measures indicate that a large portion of OPUS100 contains low-quality data. Approximately 10% of the data has COMETKIWI scores below 0.25, indicating very poor translation quality. Additionally, between 10% to 30% of the data falls within the 0.25-0.5 score range, which is still considered sub-optimal. Regarding sentence length, we find that over 20% of the OPUS100 data consists of very short sentences, with a length of no more than five tokens. For language identification, 13% to 25% of the data is removed due to incorrect language identification results in one of the two parallel sentences.

**Low-resource languages suffer more from low-quality issues.** For low-resource languages like Hindi, there are fewer high-quality parallel sentences compared to high-resource languages such as Spanish. Analysis of translation quality indicates that the English-Hindi pair has less than 20% of parallel sentences with high COMETWIKI scores (0.75-1), whereas the English-Spanish pair has around 45%. For sentence length, the English-Hindi pair contains 10% more short sentences (0-5 tokens) compared to high-resource language pairs. Moreover, both the English-Arabic and English-Hindi pairs exhibit about 10% more parallel sentences that may be in the wrong languages.

These comprehensive findings underscore the critical importance of data quality when exploiting parallel corpora for mLLM training.

### 4.2 Effect of Quality

Table 3 presents the performance of BLOOM-7B1 after continued pretraining with the machine translation objective, using 10K parallel corpora with various quality filtering strategies.

**Parallel corpora containing noisy translations still improve results.** Comparing the results of the experiment with  $\tau_c = 0$  (ID 1) to the original model (ID 0), there’s an average improvement of 0.4% for all tasks. The most notable improvements are observed in both bilingual tasks (en-xx) and in-language tasks. However, generating English for bilingual tasks yields degraded or marginally improved results. Experiment 0 exhibits 0.7% and 1.1% decrements in FLORES and MUSE respectively, with only a 0.3% improvement in MLQA.

**Filtering out noisy translations leads to notable improvements.** When  $\tau_c = 0.5$ , the average performance rises from 53.2% to 53.7%. Further refinement to  $\tau_c = 0.75$  achieves an additional 0.3% improvement. These improvements are consistently observed across all evaluated tasks. In the optimal setting (ID 5), there’s a 1.2% improvement compared to BLOOM-7B1 (ID 0). The improvements corroborate the reliability of COMETKIWI as a metric for filtering low-quality translations.

**Filtering short sentences yields slightly worse results than using unfiltered data.** The experiment with filtering short sentences (ID 3) achieves comparable or slightly worse results compared to that without filtering short sentences (ID 5). This suggests that short sentences, whether at the word or phrase level, may offer some benefits for sentence-level tasks.

| ID | MODEL     |          |     | FLORES      |             | MUSE        |             | MLQA        |             | XQUAD       | SIB         | AVG         |
|----|-----------|----------|-----|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|    |           |          |     | EN-XX       | XX-EN       | EN-XX       | XX-EN       | EN-XX       | XX-EN       |             |             |             |
| 0  | BLOOM-7B1 |          |     | 69.1        | <b>72.4</b> | 43.1        | 53.7        | 36.4        | 42.7        | 47.2        | 58.1        | 52.8        |
|    | $\tau_c$  | $\tau_l$ | LID |             |             |             |             |             |             |             |             |             |
| 1  | 0         | 0        | ✓   | 69.7        | 71.7        | 44.4        | 52.6        | 37.8        | 43.0        | 47.7        | 58.8        | 53.2        |
| 2  | 0.5       | 0        | ✓   | 69.9        | 72.1        | 45.0        | 53.0        | <b>38.1</b> | 43.7        | 48.1        | <b>59.8</b> | 53.7        |
| 3  | 0.75      | 5        | ✓   | 70.3        | 72.1        | <b>45.7</b> | 53.6        | <b>38.1</b> | 43.5        | 47.8        | 59.2        | 53.8        |
| 4  | 0.75      | 0        | ✗   | <b>70.5</b> | 72.1        | 44.9        | 53.7        | 37.7        | <b>44.0</b> | <b>48.3</b> | 59.6        | 53.9        |
| 5  | 0.75      | 0        | ✓   | 70.3        | 72.3        | 45.5        | <b>53.9</b> | 38.0        | 43.9        | <b>48.3</b> | 59.5        | <b>54.0</b> |

Table 3: Performance (%) of BLOOM-7B1 after continued pretraining with the machine translation objective using 10K parallel corpora with various quality filtering strategies. Parameters include  $\tau_c$  for COMETWIKI score threshold,  $\tau_l$  for sentence length threshold, and LID indicating the adoption of language identification filtering.

**Using data with language identification filtering results in only a 0.1% improvement on average.** A comparison of experimental outcomes with and without language identification filtering (ID 4 and 5) reveals that using data with language identification filtering yields merely a 0.1% improvement on average. The most notable performance difference is observed in the MUSE task, where using data with language identification filtering leads to improvements of 0.6% (en-xx) and 0.2% (xx-en). This marginal enhancement may be attributed to the presence of sentences in similar languages within OPUS100, which exhibit minor linguistic variations compared to the true language. These variations could potentially have a slight negative impact on word-level translations while having little impact on sentence-level tasks.

## 5 Quantity

### 5.1 Effect of Quantity Across Tasks

Based on Table 4, which shows the performance of BLOOM-7B1 after continued pretraining with the machine translation objective using different amounts of parallel sentences, we can derive the following key findings:

**Adding merely 1K parallel sentences helps.** Exploiting 1K parallel sentences for continued pretraining improves the overall average score by 1%. This increase is observed across most tasks, with notable improvements in FLORES (en-xx), MUSE (en-xx), and SIB.

**Using 10K parallel sentences leads to the optimal performance.** The best overall performance is achieved with 10K parallel sentences, resulting in an average score of 54.0%. This setting yields the highest scores in MUSE and SIB.

**More data achieves comparable results.** Increasing the number of parallel sentences beyond 10K results in comparable performance. Specifically, using 25K or 50K parallel sentences yields average scores of 53.9%, which are very close to the score obtained with 10K sentences.

The analysis suggests that continued pretraining with a moderate amount of parallel sentences (around 10K) yields the best overall improvement in performance for the BLOOM-7B1 model across various tasks.

### 5.2 Effect of Quantity Across Languages

We delve deeper into the influence of parallel corpora quantity across various languages, as depicted in Table 5.

**Using 10K parallel sentences achieves optimal performance across most languages.** For the majority of languages, except Vietnamese (vi) and Chinese (zh), the highest performance is obtained with 10K parallel sentences. Even for Vietnamese and Chinese, leveraging 10K parallel sentences can yield comparable results. These observations align with the findings in §5.1.

**Different languages exhibit varying appetites for parallel corpora.** Across most languages, increasing the number of parallel sentences used for continued pretraining generally leads to incremental improvements in performance. However, for Hindi (hi) and Chinese (zh), transitioning from 1K to 10K parallel sentences does not yield improvement. This phenomenon may be attributed to BLOOM-7B1’s limited proficiency in these languages compared to others, as reflected in the results of the original BLOOM-7B1 model ( $|SENT|=0$ ).

| SENT | FLORES      |             | MUSE        |             | MLQA        |             | XQUAD       | SIB         | AVG         |
|------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|      | EN-XX       | XX-EN       | EN-XX       | XX-EN       | EN-XX       | XX-EN       |             |             |             |
| 0    | 69.1        | <b>72.4</b> | 43.1        | 53.7        | 36.4        | 42.7        | 47.2        | 58.1        | 52.8        |
| 1K   | 70.0        | 72.2        | 45.3        | 53.6        | <b>38.2</b> | 43.6        | 47.9        | 59.2        | 53.8        |
| 5K   | 70.3        | 72.2        | 45.4        | 53.5        | <b>38.2</b> | 43.8        | 48.2        | <b>59.5</b> | 53.9        |
| 10K  | 70.3        | 72.3        | <b>45.5</b> | <b>53.9</b> | 38.0        | 43.9        | 48.3        | <b>59.5</b> | <b>54.0</b> |
| 25K  | 70.3        | 72.2        | 45.1        | 53.8        | 38.0        | <b>44.0</b> | <b>48.4</b> | <b>59.5</b> | 53.9        |
| 50K  | <b>70.4</b> | 72.2        | 45.1        | 53.8        | 38.1        | 43.7        | 48.3        | <b>59.5</b> | 53.9        |

Table 4: Performance (%) of BLOOM-7B1 after continued pretraining with the machine translation objective using varying amounts of parallel sentences, obtained with the best filtering strategy (ID 5) as shown in Table 3. |SENT| indicates the number of parallel sentences used for continued pretraining, with |SENT|=0 representing the original BLOOM-7B1 model.

|     | ar          | es          | hi          | vi          | zh          |
|-----|-------------|-------------|-------------|-------------|-------------|
| 0   | 49.5        | 57.7        | 46.5        | 63.8        | 46.7        |
| 1K  | 50.8        | 58.1        | <b>47.7</b> | 64.3        | 47.8        |
| 5K  | 51.2        | 58.2        | 47.6        | 64.6        | <b>47.9</b> |
| 10K | <b>51.3</b> | <b>58.4</b> | <b>47.7</b> | 64.6        | 47.8        |
| 25K | 51.2        | 58.2        | <b>47.7</b> | <b>64.7</b> | 47.7        |
| 50K | 51.2        | 58.2        | <b>47.7</b> | <b>64.7</b> | 47.6        |

Table 5: Performance (%) of BLOOM-7B1 after continued pretraining with the machine translation objective using varying amounts of parallel sentences, obtained with the best filtering strategy (ID 5) as shown in Table 3. |SENT| indicates the number of parallel sentences used for continued pretraining, with |SENT|=0 representing the original BLOOM-7B1 model.

## 6 Objective

We explore the effectiveness of different objectives and their combinations, with results presented in Table 6.

**BLOOM-7B1 performs well on English generation tasks.** The baseline BLOOM-7B1 model exhibits robust performance across a spectrum of evaluation tasks, notably excelling in English generation tasks such as FLORES (xx-en) and MUSE (xx-en). Further exploitation of parallel corpora fails to yield any discernible improvement.

**MT emerges as the top performer.** The MT objective consistently outperforms the baseline BLOOM-7B1 model, showcasing an average improvement of 1.2%. Moreover, MT achieves the highest performance in 5 out of 8 evaluated tasks.

**TLM exhibits limited effectiveness.** While TLM shows slight improvements on average (0.2%), primarily driven by enhancements in tasks like MUSE (en-xx), MLQA (xx-en), XQUAD, and

SIB, it also leads to degradation in tasks including FLORES and MUSE (xx-en).

**XSS achieves strong performance for classification.** Using the XSS objective improves BLOOM-7B1 by 0.7%, though it performs 0.5% worse than MT. The major decrease is observed in translation tasks, especially from English to other languages. However, XSS can still slightly improve translation tasks compared to BLOOM-7B1. Notably, XSS achieves 0.3% better performance on SIB, highlighting its effectiveness for classification.

**Combining training objectives does not provide large benefits.** While combinations of different objectives can improve BLOOM-7B1 by 0.2% to 1.0%, none surpass the performance of using the MT objective alone. The combination of MT and XSS is the best among the combinations, slightly worse than MT by 0.2%, but better than all other objectives. Notably, MT +XSS achieves the best results on SIB, and TLM +XSS yields the best results on MLQA (xx-en). These observations indicate that no single objective excels across all tasks.

## 7 Model Size

We explore the impact of parallel corpora on various sizes of BLOOM models, detailed in Table 7.

**Smaller models exhibit more pronounced improvements in FLORES.** Notably, BLOOM-1B7 demonstrates larger improvements compared to its larger counterparts in the FLORES task, where the prompt is the same as the one used during instruction tuning with the MT objective. This is attributed to the smaller models' less developed in-context learning capabilities before instruction tuning, allowing for more substantial improvements when supplemented with parallel corpora.

| MODEL        | FLORES      |             | MUSE        |             | MLQA        |             | XQUAD       | SIB         | AVG         |
|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|              | EN-XX       | XX-EN       | EN-XX       | XX-EN       | EN-XX       | XX-EN       |             |             |             |
| BLOOM-7B1    | 69.1        | <b>72.4</b> | 43.1        | 53.7        | 36.4        | 42.7        | 47.2        | 58.1        | 52.8        |
| MT           | <b>70.3</b> | 72.3        | <b>45.5</b> | <b>53.9</b> | <b>38.0</b> | 43.9        | <b>48.3</b> | 59.5        | <b>54.0</b> |
| TLM          | 67.2        | 72.2        | 44.3        | 53.0        | 36.3        | 44.4        | 47.6        | 58.7        | 53.0        |
| XSS          | 69.4        | 72.2        | 43.7        | 53.5        | 37.0        | 44.2        | <b>48.3</b> | 60.0        | 53.5        |
| MT +TLM      | 69.3        | 72.1        | 44.1        | 53.2        | 36.8        | 43.8        | 47.2        | 59.5        | 53.2        |
| MT +XSS      | 70.3        | 72.1        | 44.9        | 53.3        | 37.4        | 44.5        | 47.9        | <b>60.4</b> | 53.8        |
| TLM +XSS     | 67.7        | 72.2        | 43.0        | 52.5        | 34.9        | <b>45.6</b> | 48.2        | 60.0        | 53.0        |
| MT +TLM +XSS | 69.5        | 72.1        | 44.2        | 53.2        | 36.1        | 45.1        | 47.7        | 59.0        | 53.4        |

Table 6: Performance (%) of BLOOM-7B1 after continued pretraining with different objectives and their combinations using 10K parallel corpora, obtained with the best filtering strategy (ID 5) as shown in Table 3.

| MODEL           | FLORES      |              | MUSE        |              | MLQA        |              | XQUAD        | SIB         | AVG         |
|-----------------|-------------|--------------|-------------|--------------|-------------|--------------|--------------|-------------|-------------|
|                 | EN-XX       | XX-EN        | EN-XX       | XX-EN        | EN-XX       | XX-EN        |              |             |             |
| BLOOM-7B1       | 69.1        | 72.4         | 43.1        | 53.7         | 36.4        | 42.7         | 47.2         | 58.1        | 52.8        |
| + Parallel Data | 70.3        | 72.3         | 45.5        | 53.9         | 38.0        | 43.9         | 48.3         | 59.5        | 54.0        |
| $\Delta$        | <b>01.2</b> | <b>-00.1</b> | <b>02.4</b> | <b>00.2</b>  | <b>01.6</b> | <b>01.2</b>  | <b>01.0</b>  | <b>01.4</b> | <b>01.2</b> |
| BLOOM-3B        | 64.0        | 68.9         | 39.7        | 50.9         | 29.4        | 26.2         | 32.7         | 54.5        | 45.8        |
| + Parallel Data | 65.0        | 69.1         | 41.4        | 51.6         | 30.9        | 26.7         | 34.5         | 56.9        | 47.0        |
| $\Delta$        | <b>01.0</b> | <b>00.2</b>  | <b>01.8</b> | <b>00.7</b>  | <b>01.5</b> | <b>00.5</b>  | <b>01.8</b>  | <b>02.4</b> | <b>01.2</b> |
| BLOOM-1B7       | 59.0        | 65.8         | 37.2        | 48.5         | 20.0        | 22.2         | 24.8         | 53.0        | 41.3        |
| + Parallel Data | 61.1        | 65.7         | 38.9        | 48.0         | 20.8        | 20.9         | 24.4         | 53.0        | 41.6        |
| $\Delta$        | <b>02.0</b> | <b>-00.1</b> | <b>01.6</b> | <b>-00.6</b> | <b>00.8</b> | <b>-01.3</b> | <b>-00.3</b> | <b>00.0</b> | <b>00.3</b> |

Table 7: Effect of parallel corpora on BLOOM models of different sizes across various tasks. ‘+ Parallel Data’ indicates continued pretraining of the given mLLM with the MT objective, using 10K parallel corpora obtained with the best filtering strategy (ID 5) as shown in Table 3.

**Larger models excel in diverse tasks.** Conversely, larger models generally demonstrate greater enhancements in tasks beyond machine translation. Both BLOOM-7B1 and BLOOM-3B exhibit a 1.2% improvement compared to their original mLLMs, while BLOOM-1B7 shows a slight 0.3% improvement. Specifically, BLOOM-7B1 and BLOOM-3B display notable improvements in tasks except for FLORES, while BLOOM-1B7 achieves comparable or even worse results.

These findings demonstrate that when leveraging parallel corpora to enhance mLLMs, larger models not only exhibit improvements in direct tasks, such as machine translation, but also demonstrate a more substantial overall enhancement across a variety of tasks. In contrast, smaller models primarily show benefits in direct tasks. This difference can be attributed to the superior cross-task transferability of larger mLLMs, where insights gained from parallel corpora in one task contribute to improved performance in others.

## 8 Conclusion

This paper investigates the impact of four critical factors – data quality, data quantity, objectives, and mLLM sizes – on leveraging parallel corpora to enhance mLLMs across diverse languages and tasks. Our findings underscore the crucial importance of filtering out noisy translations to procure high-quality training data for improving mLLMs. Surprisingly, even a relatively modest dataset of 10K samples can yield promising results. Furthermore, our analysis shows that employing the machine translation objective leads to optimal outcomes. Importantly, larger models exhibit a greater capacity to benefit from parallel corpora, achieving more substantial improvements. This study provides a comprehensive recipe for effectively leveraging parallel corpora to enhance mLLMs. These insights significantly contribute to advancing the understanding and optimization of mLLMs across different languages and tasks.



## 550 Limitations

551 Due to limited computational resources, we opted  
552 not to explore full-parameter continued pretraining  
553 for leveraging parallel corpora. Instead, we focused  
554 on LoRA, drawing on insights from previous stud-  
555 ies. Additionally, our investigation is restricted to  
556 the BLOOM series, and we did not extend our anal-  
557 ysis to other mLLMs. Furthermore, we did not also  
558 explore mLLMs larger than 7B1.

## 559 References

560 David Ifeoluwa Adelani, Hannah Liu, Xiaoyu Shen,  
561 Nikita Vassilyev, Jesujoba O. Alabi, Yanke Mao, Hao-  
562 nan Gao, and En-Shiun Annie Lee. 2023. [SIB-200: A simple, inclusive, and big evaluation dataset for topic classification in 200+ languages and dialects](#). *CoRR*, abs/2309.07445.

566 Kabir Ahuja, Rishav Hada, Millicent Ochieng, Prachi  
567 Jain, Harshita Diddee, Samuel Maina, Tanuja Ganu,  
568 Sameer Segal, Maxamed Axmed, Kalika Bali, and  
569 Sunayana Sitaram. 2023. [MEGA: multilingual evaluation of generative AI](#). *CoRR*, abs/2303.12528.

571 Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pe-  
572 dro Henrique Martins, João Alves, M. Amin Farajian,  
573 Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta  
574 Agrawal, Pierre Colombo, José G. C. de Souza, and  
575 André F. T. Martins. 2024. [Tower: An open multi-lingual large language model for translation-related tasks](#). *CoRR*, abs/2402.17733.

578 Mikel Artetxe, Itziar Aldabe, Rodrigo Agerri, Olatz  
579 Perez-de-Viñaspre, and Aitor Soroa. 2022. [Does corpus quality really matter for low-resource languages?](#) In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11*, pages 7383–7390. Association for Computational Linguistics.

586 Mikel Artetxe, Sebastian Ruder, and Dani Yogatama.  
587 2020. [On the cross-lingual transferability of mono-lingual representations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 4623–4637. Association for Computational Linguistics.

593 Samuel Cahyawijaya, Holy Lovenia, Tiezheng Yu,  
594 Willy Chung, and Pascale Fung. 2023. [Instruct-align: Teaching novel languages with to llms through alignment-based cross-lingual instruction](#). *CoRR*, abs/2305.13627.

598 Hao Chen, Yiming Zhang, Qi Zhang, Hantao Yang, Xi-  
599 aomeng Hu, Xuetao Ma, Yifan Yanggong, and Junbo  
600 Zhao. 2023. [Maybe only 0.5% data is needed: A preliminary exploration of low training data instruction tuning](#). *CoRR*, abs/2305.09246.

Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham  
Singhal, Wenhui Wang, Xia Song, Xian-Ling Mao,  
Heyan Huang, and Ming Zhou. 2021a. [Infoxlm: An information-theoretic framework for cross-lingual language model pre-training](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 3576–3588. Association for Computational Linguistics.

Zewen Chi, Li Dong, Bo Zheng, Shaohan Huang, Xian-  
Ling Mao, Heyan Huang, and Furu Wei. 2021b. [Improving pretrained cross-lingual language models via self-labeled word alignment](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 3418–3430. Association for Computational Linguistics.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal,  
Vishrav Chaudhary, Guillaume Wenzek, Francisco  
Guzmán, Edouard Grave, Myle Ott, Luke Zettle-  
moyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 8440–8451. Association for Computational Linguistics.

Alexis Conneau and Guillaume Lample. 2019. [Cross-lingual language model pretraining](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 7057–7067.

Marta R. Costa-jussà, James Cross, Onur Çelebi,  
Maha Elbayad, Kenneth Heafield, Kevin Heffernan,  
Elahe Kalbassi, Janice Lam, Daniel Licht, Jean  
Maillard, Anna Sun, Skyler Wang, Guillaume  
Wenzek, Al Youngblood, Bapi Akula, Loïc Bar-  
rault, Gabriel Mejia Gonzalez, Prangthip Hansanti,  
John Hoffman, Semarley Jarrett, Kaushik Ram  
Sadagopan, Dirk Rowe, Shannon Spruit, Chau  
Tran, Pierre Andrews, Necip Fazil Ayan, Shruti  
Bhosale, Sergey Edunov, Angela Fan, Cynthia  
Gao, Vedanuj Goswami, Francisco Guzmán, Philipp  
Koehn, Alexandre Mourachko, Christophe Rop-  
pers, Safiyyah Saleem, Holger Schwenk, and Jeff  
Wang. 2022. [No language left behind: Scaling human-centered machine translation](#). *CoRR*, abs/2207.04672.

Himanshu Gupta, Saurabh Arjun Sawant, Swaroop  
Mishra, Mutsumi Nakamura, Arindam Mitra, San-  
tosh Mashetty, and Chitta Baral. 2023. [Instruction tuned models are quick learners](#). *CoRR*, abs/2306.05539.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan  
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and  
Weizhu Chen. 2022. [Lora: Low-rank adaptation of](#)

|     |   |   |   |
|-----|---|---|---|
| 662 | <a href="#">large language models</a> . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.   |   |   |
| 663 |   |   |   |
| 664 |   |   |   |
| 665 | Junjie Hu, Melvin Johnson, Orhan Firat, Aditya Siddhant, and Graham Neubig. 2021. <a href="#">Explicit alignment objectives for multilingual bidirectional encoders</a> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021</i> , pages 3633–3643. Association for Computational Linguistics.   |   |   |
| 666 |   |   |   |
| 667 |   |   |   |
| 668 |   |   |   |
| 669 |   |   |   |
| 670 |   |   |   |
| 671 |   |   |   |
| 672 |   |   |   |
| 673 |   |   |   |
| 674 | Haoyang Huang, Yaobo Liang, Nan Duan, Ming Gong, Linjun Shou, Daxin Jiang, and Ming Zhou. 2019. <a href="#">Unicoder: A universal language encoder by pre-training with multiple cross-lingual tasks</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 2019, Hong Kong, China, November 3-7, 2019</i> , pages 2485–2494. Association for Computational Linguistics.  |   |   |
| 675 |   |   |   |
| 676 |   |   |   |
| 677 |   |   |   |
| 678 |   |   |   |
| 679 |   |   |   |
| 680 |   |   |   |
| 681 |   |   |   |
| 682 |   |   |   |
| 683 |   |   |   |
| 684 | Ayyoob Imani, Peiqin Lin, Amir Hossein Kargaran, Silvia Severini, Masoud Jalili Sabet, Nora Kassner, Chunlan Ma, Helmut Schmid, André F. T. Martins, François Yvon, and Hinrich Schütze. 2023. <a href="#">Glot500: Scaling multilingual corpora and language models to 500 languages</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 1082–1117. Association for Computational Linguistics.   |   |   |
| 685 |   |   |   |
| 686 |   |   |   |
| 687 |   |   |   |
| 688 |   |   |   |
| 689 |   |   |   |
| 690 |   |   |   |
| 691 |   |   |   |
| 692 |   |   |   |
| 693 |   |   |   |
| 694 | Amir Hossein Kargaran, Ayyoob Imani, François Yvon, and Hinrich Schütze. 2023. <a href="#">Glotlid: Language identification for low-resource languages</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 6155–6218. Association for Computational Linguistics.  |   |   |
| 695 |   |   |   |
| 696 |   |   |   |
| 697 |   |   |   |
| 698 |   |   |   |
| 699 |   |   |   |
| 700 |   |   |   |
| 701 | Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Javier Ortiz Suárez, Iroko Orife, Kelechi Ogejeji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Balli, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2022. <a href="#">Quality at a glance: An audit of web-crawled multilingual</a> |   |   |
| 702 |   |   |   |
| 703 |   |   |   |
| 704 |   |   |   |
| 705 |   |   |   |
| 706 |   |   |   |
| 707 |   |   |   |
| 708 |   |   |   |
| 709 |   |   |   |
| 710 |   |   |   |
| 711 |   |   |   |
| 712 |   |   |   |
| 713 |   |   |   |
| 714 |   |   |   |
| 715 |   |   |   |
| 716 |   |   |   |
| 717 |   |   |   |
| 718 |   |   |   |
| 719 |   |   |   |
| 720 |   |   |   |
|     |   | <a href="#">datasets</a> . <i>Trans. Assoc. Comput. Linguistics</i> , 10:50–72. | 721<br>722                                    |
|     | Guillaume Lample, Alexis Conneau, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018. <a href="#">Word translation without parallel data</a> . In <i>6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings</i> . OpenReview.net.   |   | 723<br>724<br>725<br>726<br>727<br>728<br>729 |
|     | Patrick S. H. Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2020. <a href="#">MLQA: evaluating cross-lingual extractive question answering</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 7315–7330. Association for Computational Linguistics.  |   | 730<br>731<br>732<br>733<br>734<br>735<br>736 |
|     | Chong Li, Shaonan Wang, Jiajun Zhang, and Chengqing Zong. 2023. <a href="#">Align after pre-train: Improving multilingual generative models with cross-lingual alignment</a> . <i>CoRR</i> , abs/2311.08089.  |   | 737<br>738<br>739<br>740                      |
|     | Peiqin Lin, Chengzhi Hu, Zheyu Zhang, André F. T. Martins, and Hinrich Schütze. 2024a. <a href="#">mplm-sim: Better cross-lingual similarity and transfer in multilingual pretrained language models</a> . In <i>Findings of the Association for Computational Linguistics: EACL 2024, St. Julian’s, Malta, March 17-22, 2024</i> , pages 276–310. Association for Computational Linguistics.   |   | 741<br>742<br>743<br>744<br>745<br>746<br>747 |
|     | Peiqin Lin, Shaoxiong Ji, Jörg Tiedemann, André F. T. Martins, and Hinrich Schütze. 2024b. <a href="#">Mala-500: Massive language adaptation of large language models</a> . <i>CoRR</i> , abs/2401.13303.   |   | 748<br>749<br>750<br>751                      |
|     | Peiqin Lin, André F. T. Martins, and Hinrich Schütze. 2024c. <a href="#">Xampler: Learning to retrieve cross-lingual in-context examples</a> . <i>Preprint</i> , arXiv:2405.05116.  |   | 752<br>753<br>754                             |
|     | Yihong Liu, Peiqin Lin, Mingyang Wang, and Hinrich Schütze. 2023. <a href="#">OFA: A framework of initializing unseen subword embeddings for efficient large-scale multilingual continued pretraining</a> . <i>CoRR</i> , abs/2311.08849.   |   | 755<br>756<br>757<br>758<br>759               |
|     | Sheng Lu, Irina Bigoulaeva, Rachneet Sachdeva, Harish Tayyar Madabushi, and Iryna Gurevych. 2023. <a href="#">Are emergent abilities in large language models just in-context learning?</a> <i>CoRR</i> , abs/2309.01809.   |   | 760<br>761<br>762<br>763                      |
|     | Xuan Ouyang, Shuohuan Wang, Chao Pang, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. 2020. <a href="#">ERNIE-M: enhanced multilingual representation by aligning cross-lingual semantics with monolingual corpora</a> . <i>CoRR</i> , abs/2012.15674.   |   | 764<br>765<br>766<br>767<br>768               |
|     | Laura Cabello Piqueras and Anders Søgaard. 2022. <a href="#">Are pretrained multilingual models equally fair across languages?</a> In <i>Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022</i> , pages 3597–3605. International Committee on Computational Linguistics.  |   | 769<br>770<br>771<br>772<br>773<br>774<br>775 |

|     |  |     |
|-----|--|-----|
| 776 | Surangika Ranathunga, Nisansa de Silva, Menan Velayuthan, Aloka Fernando, and Charitha Rathnayake. 2024. <a href="#">Quality does matter: A detailed look at the quality and utility of web-mined parallel corpora</a> . <i>CoRR</i> , abs/2402.07446.   | 836 |
| 777 |  | 837 |
| 778 |  | 838 |
| 779 |  | 839 |
| 780 |  | 840 |
| 781 | Ricardo Rei, Marcos V. Treviso, Nuno Miguel Guerreiro, Chrysoula Zerva, Ana C. Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte M. Alves, Luísa Coheur, Alon Lavie, and André F. T. Martins. 2022. <a href="#">Cometkiwi: Ist-unbabel 2022 submission for the quality estimation shared task</a> . In <i>Proceedings of the Seventh Conference on Machine Translation, WMT 2022, Abu Dhabi, United Arab Emirates (Hybrid), December 7-8, 2022</i> , pages 634–645. Association for Computational Linguistics.   | 841 |
| 782 |  | 842 |
| 783 |  | 843 |
| 784 |  | 844 |
| 785 |  | 845 |
| 786 |  | 846 |
| 787 |  | 847 |
| 788 |  | 848 |
| 789 |  | 849 |
| 790 |  | 850 |
| 791 | Machel Reid and Mikel Artetxe. 2022a. <a href="#">On the role of parallel data in cross-lingual transfer learning</a> . <i>CoRR</i> , abs/2212.10173.  | 851 |
| 792 |  | 852 |
| 793 |  | 853 |
| 794 | Machel Reid and Mikel Artetxe. 2022b. <a href="#">PARADISE: exploiting parallel data for multilingual sequence-to-sequence pretraining</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022</i> , pages 800–810. Association for Computational Linguistics.  | 854 |
| 795 |  | 855 |
| 796 |  | 856 |
| 797 |  | 857 |
| 798 |  | 858 |
| 799 |  | 859 |
| 800 |  | 860 |
| 801 |  | 861 |
| 802 | Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilic, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Amanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, and et al. 2022. <a href="#">BLOOM: A 176b-parameter open-access multilingual language model</a> . <i>CoRR</i> , abs/2211.05100. | 862 |
| 803 |  | 863 |
| 804 |  | 864 |
| 805 |  | 865 |
| 806 |  | 866 |
| 807 |  | 867 |
| 808 |  | 868 |
| 809 |  | 869 |
| 810 |  | 870 |
| 811 |  | 871 |
| 812 |  | 872 |
| 813 |  | 873 |
| 814 |  | 874 |
| 815 |  | 875 |
| 816 |  | 876 |
| 817 |  | 877 |
| 818 |  | 878 |
| 819 |  | 879 |
| 820 |  | 880 |
| 821 | Manli Shu, Jiong Xiao Wang, Chen Zhu, Jonas Geiping, Chaowei Xiao, and Tom Goldstein. 2023. <a href="#">On the exploitability of instruction tuning</a> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .  | 881 |
| 822 |  | 882 |
| 823 |  | 883 |
| 824 |  | 884 |
| 825 |  | 885 |
| 826 |  | 886 |
| 827 |  | 887 |
| 828 | Ahmet Üstün, Viraat Aryabumi, Zheng Xin Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blunsom, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. 2024. <a href="#">Aya model: An instruction finetuned open-access multilingual language model</a> . <i>CoRR</i> , abs/2402.07827.  | 888 |
| 829 |  | 889 |
| 830 |  | 890 |
| 831 |  |     |
| 832 |  |     |
| 833 |  |     |
| 834 |  |     |
| 835 |  |     |
|     | Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023. <a href="#">How far can camels go? exploring the state of instruction tuning on open resources</a> . <i>CoRR</i> , abs/2306.04751.  |     |
|     | Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. <a href="#">Emergent abilities of large language models</a> . <i>Trans. Mach. Learn. Res.</i> , 2022.   |     |
|     | Jerry W. Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. 2023. <a href="#">Larger language models do in-context learning differently</a> . <i>CoRR</i> , abs/2303.03846.   |     |
|     | Xiangpeng Wei, Rongxiang Weng, Yue Hu, Luxi Xing, Heng Yu, and Weihua Luo. 2021. <a href="#">On learning universal representations across languages</a> . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.   |     |
|     | Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. <a href="#">A paradigm shift in machine translation: Boosting translation performance of large language models</a> . <i>CoRR</i> , abs/2309.11674.   |     |
|     | Jian Yang, Shuming Ma, Dongdong Zhang, Shuangzhi Wu, Zhoujun Li, and Ming Zhou. 2020. <a href="#">Alternating language modeling for cross-lingual pre-training</a> . In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020</i> , pages 9386–9393. AAAI Press.  |     |
|     | Biao Zhang, Philip Williams, Ivan Titov, and Rico Senrich. 2020. <a href="#">Improving massively multilingual neural machine translation and zero-shot translation</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 1628–1639. Association for Computational Linguistics.  |     |
|     | Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. <a href="#">LIMA: less is more for alignment</a> . <i>CoRR</i> , abs/2305.11206.   |     |
|     | Wenhao Zhu, Yunzhe Lv, Qingxiu Dong, Fei Yuan, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2023. <a href="#">Extrapolating large language models to non-english by aligning languages</a> . <i>CoRR</i> , abs/2308.04948.  |     |

891 **A Prompt**

892 The prompts of FLORES, MUSE, MLQA,  
893 XQuAD, and SIB are shown as follows:

894 **FLORES/MUSE**

895 Translate the following  
896 text from [SOURCE\_LANG]  
897 to [TARGET\_LANG].\nText:  
898 [SOURCE\_TEXT]\nTranslation:  
899 [TARGET\_TEXT]

900 **MLQA/XQuAD**

901 [Passage] \nQ:  
902 [Question]\nA: [Answer]

903 **SIB**

904 The topic of the news  
905 [Passage] is [Label]