# *LangGPS*: Language Separability Guided Data Pre-Selection for Joint Multilingual Instruction Tuning

**Yangfan Ye[1], Xiaocheng Feng[1,2*], Xiachong Feng[3*], Lei Huang[1], Weitao Ma[1], Qichen Hong[4], Yunfei Lu[4], Duyu Tang[4], Dandan Tu[4*], Bing Qin[1,2]**

[1]Harbin Institute of Technology, [2]Peng Cheng Laboratory, [3]The University of Hong Kong, [4]Huawei Technologies Co., Ltd
{yfye,xcfeng,lhuang,wtma,qinb}@ir.hit.edu.cn, fengxc@hku.hk, {hongqichen,luyunfei6,tangduyu,tudandan}@huawei.com

## Abstract

Joint multilingual instruction tuning is a widely adopted approach to improve the multilingual instruction-following ability and downstream performance of large language models (LLMs), but the resulting multilingual capability remains highly sensitive to the composition and selection of the training data. Existing selection methods, often based on features like text quality, diversity, or task relevance, typically overlook the intrinsic linguistic structure of multilingual data. In this paper, we propose *LangGPS*, a lightweight two-stage pre-selection framework guided by *language separability*—a signal that quantifies how well samples in different languages can be distinguished in the model's representation space. *LangGPS* first filters training data based on separability scores and then refines the subset using existing selection methods. Extensive experiments across six benchmarks and 22 languages demonstrate that applying *LangGPS* on top of existing selection methods improves their effectiveness and generalizability in multilingual training, especially for understanding tasks and low-resource languages. Further analysis reveals that highly separable samples facilitate the formation of clearer language boundaries and support faster adaptation, while low-separability samples tend to function as bridges for cross-lingual alignment. Besides, we also find that language separability can serve as an effective signal for multilingual curriculum learning, where interleaving samples with diverse separability levels yields stable and generalizable gains. Together, we hope our work offers a new perspective on data utility in multilingual contexts and support the development of more linguistically informed LLMs.

- *Code:* https://github.com/YYF-Tommy/LangGPS

## 1 Introduction

Large language models (LLMs) have demonstrated remarkable generalization and adaptability across a wide range of downstream tasks in multilingual scenarios (Ye, Tao, and Kong 2023; Qin et al. 2024; Zhang et al. 2023; Ye et al. 2024a,b). A central enabler of such capabilities is joint multilingual instruction tuning, in which LLMs are trained on instruction-following data covering multiple languages (Ouyang et al. 2022; Shaham et al. 2024; Huo et al.

2025; Ye et al. 2025). This paradigm allows models to understand and follow instructions in multilingual scenarios, fostering more inclusive access to language technologies.

Recent work has emphasized the importance of multilingual instruction datasets that encompass diverse languages, domains, and task formats (Muennighoff et al. 2022; Zhu et al. 2023; Li et al. 2023; Singh et al. 2024). These efforts have laid a solid foundation for developing LLMs that can effectively align with user intent across linguistic boundaries. Despite these advances, the success of multilingual instruction tuning remains highly sensitive to the composition and selection of the training data. Consequently, **Data Selection** has emerged as a critical strategy for balancing effectiveness and efficiency in large-scale model training. Existing approaches can be typically grouped into two main categories: feature-based methods, which prioritize training samples based on characteristics such as data quality and diversity (Li et al. 2024; Bukharin et al. 2024); and target-dependent methods, which select data according to its estimated relevance to specific target tasks (Xie et al. 2023; Xia et al. 2024). However, little attention has been paid to the intrinsic linguistic structure inherent in multilingual data.

In multilingual scenarios, where the model must learn to follow instructions across diverse and often low-resource languages, we argue that a key prerequisite for effective multilingual instruction following is the model's ability to form and maintain clear linguistic boundaries. Without sufficiently distinguishing between languages, the model risks conflating linguistic patterns, especially in typical low-resource languages. To formalize this intuition, we introduce the notion of language separability, which quantifies the extent to which the data samples in different languages are distinguishable in the model's representation space. Highly separable samples tend to exhibit well-structured, language-specific representations, whereas poorly separable ones often entangle with other languages in the representation space, lacking clear linguistic boundaries. As shown in Figure 1 (a), our preliminary experiments (see details in Appendix A.1) on the multilingual MMLU dataset reveal that, under limited training data ($< 2000$), training with highly separable samples leads to significantly better performance compared to low-separability or randomly selected samples. However, as Figure 1 (b) illustrates, highly separable samples may also suffer from over-similarity (excessive re-
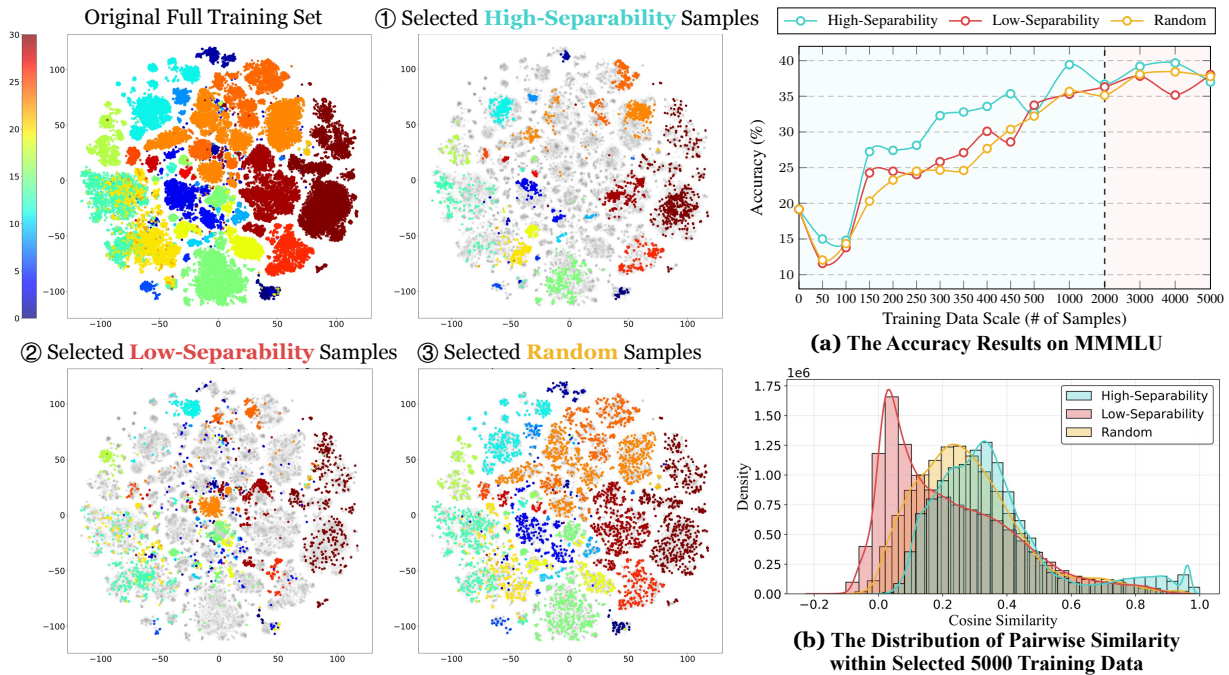
*Corresponding Authors

Figure 1: The left shows the t-SNE visualization of samples selected by different strategies under the same language distribution; each color denotes one of 31 languages. The representations are last-token hidden states from *LLaMA-3.1-8B*. On the right, panel (a) presents the performance of *LLaMA-3.1-8B* after training with the three types of samples on *MMMLU* across training sizes; while panel (b) illustrates the pairwise similarity distributions of selected data. **(See detailed settings in Appendix A.1)**

semblance among samples), introducing a lack of diversity, which is another critical factor in data selection.

Building on these insights, we propose *LangGPS*, a **Lang**uage separability **G**uided data **P**re-**S**election strategy for multilingual instruction tuning. *LangGPS* adopts a two-stage framework: (1) a lightweight pre-selection stage that filters training data based on estimated language separability, and (2) a fine selection stage that refines the selected subset using existing data selection methods, such as diversity-based, task relevance-based, or gradient-based approaches. This design integrates scalability, linguistic structure awareness, and full compatibility with existing selection strategies that focus on different aspects, offering a more principled and flexible framework for multilingual instruction tuning.

We conduct extensive experiments across six benchmarks covering both understanding and generation tasks, spanning 22 evaluated languages and two representative LLMs. Our results find that existing data selection methods often overlook the intrinsic linguistic structure inherent in multilingual training data, leading to inconsistent gains in joint multilingual training. By leveraging language separability as a guiding signal, *LangGPS* can be seamlessly integrated with existing selection methods and improves their effectiveness and generalizability, especially for understanding tasks and low-resource languages. Further analysis reveals that highly separable samples are critical for language-specific modeling and fast adaptation, helping models form and main clearer language boundaries, while low-separability samples may serve as valuable bridges for cross-lingual alignment

due to their entangled representations across languages. Moreover, we demonstrate that language separability not only informs data selection, but also enables effective curriculum learning: interleaving training samples with diverse separability levels allows models to absorb both structured and entangled linguistic signals, improving training effectiveness. These findings highlight the importance of considering linguistic structure in multilingual training pipelines, and we hope they inspire future work on more linguistically aware approaches in multilingual community.

## 2   Related Work

**Joint Multilingual Instruction Fine-Tuning.** While early work shows that fine-tuning on English-only data can yield cross-lingual generalization (Muennighoff et al. 2022), multilingual instruction tuning—i.e., fine-tuning on data from multiple languages—has emerged as a more widely adopted strategy for improving multilingual performance (Ouyang et al. 2022). Recent studies often incorporate data augmentation or distillation techniques to enhance multilingual capabilities. For example, Li et al. (2023) addressed the issue of "translationese" by generating multilingual responses using Google Translate and ChatGPT. Some other approaches incorporate multilingual examples into English-centric tuning (Shaham et al. 2024), apply translation-based finetuning to improve semantic alignment (Ranaldi, Pucci, and Freitas 2024), combine translation data, cross-lingual tasks, and scaling laws (Zhu et al. 2023), or leverage self-distillation from

high-resource languages (Zhang et al. 2024). Beyond data-level interventions, Ye et al. (2025) explored latent-level cross-lingual interactions via a cross-lingual connection mechanism. Empirical results consistently show that incorporating multilingual examples during instruction tuning leads to significantly stronger multilingual generalization compared to purely monolingual or English-centric approaches (Shaham et al. 2024; Chen et al. 2024).

**Data Selection.** A central goal of data selection in model fine-tuning is to minimize the amount of training data while maintaining—or even improving—model performance. Early work primarily focused on ensuring data quality and diversity. For example, Touvron et al. (2023) and Zhou et al. (2023) employed elaborate filtering pipelines to emphasize the role of high-quality supervision. Data diversity has also gained increasing attention, as it is crucial for training robust and generalizable models (Bukharin et al. 2024). Lu et al. (2023) proposed measuring diversity via intention tags of instructions, while Yang et al. (2025) achieved diversity-driven data selection through sparse autoencoder. Other selection criteria include features such as sequence length (Zhao et al. 2024), task complexity (Xu et al. 2024), etc. Beyond data-intrinsic properties, a growing line of work adopts model-aware selection, identifying samples that align closely with the target domain based on gradient similarity or other interaction signals between training and target sets (Xia et al. 2024; Zhao et al. 2025).

While these approaches have advanced instruction tuning through data-level and model-level interventions, little attention has been paid to the intrinsic linguistic structure present in multilingual training data. Our work addresses this gap by investigating language separability as a guiding signal for multilingual data selection.

## 3 Measurement of Language Separability

To quantify how well the training samples from different languages are separated in representation space, we use the silhouette score (Rousseeuw 1987). Originally designed to assess clustering quality, this metric favors compact intra-group and distant inter-group structures, making it well-suited for evaluating language separability.

Let $D = \{d_1, d_2, \ldots, d_L\}$ denote a multilingual supervised instruction dataset, where each $d_l = \{p_i^l\} = \{(x_i^l, y_i^l)\}$ represents a language-specific subset corresponding to language $l \in \{1, \ldots, L\}$. Here, each sample pair $p_i^l$ consists of an input instruction $x_i^l$ and its corresponding ground-truth response $y_i^l$. Given a model, we obtain the representation for each data point by formatting the instruction–response pair $(x_i^l, y_i^l)$ using the model's training template, feeding it into the model, and extracting the hidden state of the last token. All representations are then grouped by language labels.

For each data point $p_i^l$, we compute its average distance $a(p_i^l)$ to all other points in the same language cluster $d_l$:

$$a(p_i^l) = \frac{1}{|d_l| - 1} \sum_{p_j^l \in d_l, j \neq i} \text{dist}(p_i^l, p_j^l) \qquad (1)$$

where $|d_l|$ denotes the number of data points in $d_l$, and $\text{dist}(p_i^l, p_j^l)$ represents the Euclidean distance between $p_i^l$

and $p_j^l$ in the their representation space. The value of $a(p_i^l)$ quantifies how well the sample $p_i^l$ fits within its language cluster $d_l$ (the smaller the value, the better the alignment).

Next, we measure the dissimilarity between $p_i^l$ and other clusters by computing its average distance to the samples in its nearest neighboring cluster:

$$b(p_i^l) = \min_{m \neq l} \frac{1}{|d_m|} \sum_{p_j^m \in d_m} \text{dist}(p_i^l, p_j^m) \qquad (2)$$

A larger $b(p_i^l)$ indicates that $p_i^l$ is more distinct from samples in other language clusters.

The silhouette score for each data point $p_i^l$ jointly considers intra-cluster compactness $a(p_i^l)$ and inter-cluster separation $b(p_i^l)$, and is defined as:

$$s(p_i^l) = \frac{b(p_i^l) - a(p_i^l)}{\max\{a(p_i^l), b(p_i^l)\}} \qquad (3)$$

The silhouette score ranges from $-1$ to $1$, with higher values indicating that the sample $p_i^l$ is well-clustered and clearly separated from other languages in the representation space. Our language separability guided data pre-selection strategy prioritizes the top $\rho\%$ of samples with the highest silhouette scores within each language cluster, encouraging the model to form and maintain clearer linguistic boundaries by training on selected highly separable data.

## 4 Experiments

### 4.1 Setups

**Models.** We selected two representative LLMs for our main experiments: (1) *LLaMA-3.1-8B* (Dubey et al. 2024) and (2) *Qwen2.5-7B* (Yang et al. 2024).

**Training Corpus.** We totally select 97,696 multilingual instruction pairs from *aya dataset* (Singh et al. 2024) as our full training corpus and the training corpus covers 31 languages, ensuring extensive multilingual coverage (see detailed statistics in Appendix A.2). Our training processes are conducted on *8 * A100-80GB* GPUs with the following settings: *batch size=16*, *epochs=3*, *learning rate=1.0e-5*, *warmup ratio=0.1*, and *bf16=true*. The implementation is based on *LLaMA-Factory* (Zheng et al. 2024).

**Evaluation Datasets.** We conduct experiments on 6 benchmarks, which can be categorized into:

- **Multilingual Understanding:** (1) *XNLI* (Conneau et al. 2018), a multilingual natural language inference (NLI) dataset, (2) *XStoryCloze* (Lin et al. 2022), a multilingual commonsense reasoning dataset for evaluating story understanding and (3) *MMMLU*, the multilingual version of *MMLU* (Hendrycks et al. 2020), designed to evaluate models' general knowledge.
- **Multilingual Generation:** (1) *MKQA* (Longpre, Lu, and Daiber 2021), an open-domain multilingual question answering evaluation dataset, (2) *XQuAD* (Artetxe, Ruder, and Yogatama 2020), a question answering dataset and (3) *XLSum* (Hasan et al. 2021), a multilingual abstractive summarization benchmark comprising professionally annotated article-summary pairs.

**Model: LLaMA-3.1-8B**

| Method | Multilingual Understanding | | | | | | | | | | | | Multilingual Generation | | | | | | | | | | | | Δ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | XNLI | | | | XStoryCloze | | | | MMMLU | | | | MKQA | | | | XQuAD | | | | XLSum | | | | |
| | 1% | 3% | 5% | Avg. | 1% | 3% | 5% | Avg. | 1% | 3% | 5% | Avg. | 1% | 3% | 5% | Avg. | 1% | 3% | 5% | Avg. | 1% | 3% | 5% | Avg. | |
| *Full* | | 36.6 | | | | 64.2 | | | | 37.7 | | | | 15.1 | | | | 60.8 | | | | 23.3 | | | |
| *Rand* | 25.1 | 33.0 | 35.4 | 31.2 | 43.7 | 61.4 | 61.1 | 55.4 | 34.2 | 37.4 | 38.1 | 36.5 | 14.6 | 14.4 | 14.4 | 14.5 | 63.6 | 59.9 | 59.6 | 61.1 | 21.9 | 22.8 | 22.8 | 22.5 | −6.44% |
| *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC** | 24.3 | 35.5 | 34.1 | 31.3 | 44.1 | 63.7 | 65.6 | 57.8 | 38.1 | 39.6 | 37.2 | 38.3 | 15.1 | 14.8 | 15.0 | 15.0 | 59.7 | 58.9 | 59.3 | 59.3 | 21.0 | 23.1 | 22.4 | 22.2 | −5.18% |
| *MTLD* | 28.8 | 34.9 | 34.4 | 32.7 | 34.3 | 52.9 | 47.1 | 44.8 | 25.1 | 32.5 | 29.9 | 29.1 | 15.7 | 15.9 | 14.5 | 15.4 | 60.3 | 62.5 | 57.2 | 60.0 | 19.5 | 21.9 | 24.2 | 21.8 | −11.55% |
| *Nat* | 20.4 | 40.1 | 36.2 | 32.2 | 13.3 | 40.0 | 54.8 | 36.0 | 24.5 | 35.3 | 37.6 | 32.4 | 18.5 | 17.3 | 17.1 | 17.6 | 59.3 | 61.3 | 59.5 | 60.0 | 19.6 | 20.9 | 22.1 | 20.8 | −10.80% |
| *Coh** | 28.3 | 37.1 | 36.3 | 33.9 | 42.1 | 54.5 | 51.5 | 49.4 | 34.6 | 35.9 | 36.4 | 35.6 | 18.2 | 18.1 | 18.4 | 18.2 | 64.1 | 64.6 | 64.8 | 64.5 | 22.4 | 24.0 | 23.8 | 23.4 | −1.47% |
| *Und* | 23.5 | 38.1 | 35.0 | 32.2 | 25.3 | 47.5 | 59.6 | 44.1 | 28.8 | 35.0 | 38.4 | 34.1 | 18.2 | 16.9 | 17.4 | 17.5 | 60.6 | 60.6 | 59.5 | 60.2 | 20.6 | 21.2 | 22.2 | 21.3 | −7.70% |
| *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR* | 29.5 | 27.7 | 36.9 | 31.4 | 36.7 | 39.7 | 49.5 | 42.0 | 27.3 | 30.1 | 35.1 | 30.8 | 17.7 | 16.9 | 16.4 | 17.0 | 64.7 | 63.4 | 62.4 | 63.5 | 22.1 | 23.1 | 23.5 | 22.9 | −8.66% |
| *LESS** | 33.9 | 33.7 | 37.8 | 35.1 | 65.5 | 71.2 | 68.3 | 68.3 | 40.4 | 40.6 | 40.8 | 40.6 | 18.2 | 18.0 | 17.2 | 17.8 | 65.2 | 65.6 | 64.8 | 65.2 | 21.2 | 21.7 | 22.8 | 21.9 | +4.90% |
| *Applying LangGPS (to random selection and three best well-performing baselines marked with *)* | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 26.4 | 34.7 | 36.2 | 32.5 (+1.3) | 48.7 | 72.5 | 72.5 | 64.6 (+9.2) | 33.7 | 38.2 | 39.7 | 37.2 (+0.8) | 15.4 | 15.0 | 14.7 | 15.0 (+0.5) | 62.6 | 60.2 | 62.0 | 61.6 (+0.5) | 21.9 | 22.8 | 23.8 | 22.9 (+0.4) | −2.21% |
| *KMC* | 26.1 | 38.0 | 37.6 | 33.9 (+2.6) | 42.4 | 66.3 | 69.2 | 59.3 (+1.5) | 38.3 | 40.2 | 38.8 | 39.1 (+0.7) | 15.3 | 15.1 | 13.8 | 14.7 (-0.3) | 61.8 | 60.3 | 60.6 | 60.9 (+1.6) | 21.5 | 22.2 | 23.8 | 22.5 (+0.3) | −2.86% |
| *Coh* | 33.7 | 32.8 | 38.1 | 34.9 (+1.0) | 47.1 | 47.2 | 60.8 | 51.7 (+2.3) | 37.6 | 35.4 | 37.5 | 36.9 (+1.3) | 19.1 | 18.2 | 17.6 | **18.3** (+0.1) | 63.4 | 62.6 | 63.7 | 63.2 (-1.3) | 22.2 | 24.3 | 24.3 | **23.6** (+0.2) | −0.01% |
| *LESS* | 35.6 | 37.8 | 37.7 | **37.0** (+1.9) | 66.1 | 76.8 | 74.0 | **72.3** (+4.0) | 39.9 | 40.0 | 42.2 | **40.7** (+0.1) | 17.1 | 16.5 | 17.0 | 16.9 (-0.9) | 65.6 | 64.3 | 64.3 | **64.7** (-0.5) | 22.5 | 22.8 | 23.4 | 22.9 (+1.0) | **+6.37%** |

**Model: Qwen2.5-7B**

| Method | Multilingual Understanding | | | | | | | | | | | | Multilingual Generation | | | | | | | | | | | | Δ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | XNLI | | | | XStoryCloze | | | | MMMLU | | | | MKQA | | | | XQuAD | | | | XLSum | | | | |
| | 1% | 3% | 5% | Avg. | 1% | 3% | 5% | Avg. | 1% | 3% | 5% | Avg. | 1% | 3% | 5% | Avg. | 1% | 3% | 5% | Avg. | 1% | 3% | 5% | Avg. | |
| *Full* | | 50.8 | | | | 64.1 | | | | 43.2 | | | | 14.3 | | | | 62.0 | | | | 21.8 | | | |
| *Rand* | 45.6 | 53.0 | 54.7 | 51.1 | 54.2 | 76.0 | 72.7 | 67.7 | 49.3 | 45.6 | 44.1 | 46.4 | 15.5 | 15.1 | 14.5 | 15.0 | 65.4 | 64.0 | 63.9 | 64.4 | 22.1 | 22.2 | 22.7 | 22.3 | +4.18% |
| *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC** | 51.1 | 49.8 | 51.3 | 50.8 | 74.3 | 64.4 | 73.0 | 70.6 | 43.9 | 43.6 | 46.8 | 44.8 | 13.1 | 14.2 | 13.4 | 13.6 | 61.2 | 62.0 | 62.4 | 61.9 | 22.5 | 22.0 | 21.8 | 21.9 | +1.52% |
| *MTLD** | 50.1 | 53.9 | 51.5 | 51.8 | 64.2 | 73.6 | 72.0 | 69.9 | 46.5 | 46.1 | 43.7 | 45.4 | 15.2 | 14.2 | 14.6 | 14.7 | 65.5 | 63.5 | 62.8 | 63.9 | 21.7 | 22.6 | 22.4 | 22.2 | +4.03% |
| *Nat* | 37.9 | 42.4 | 49.1 | 43.1 | 56.7 | 51.3 | 58.8 | 55.6 | 46.5 | 44.1 | 42.5 | 44.4 | 15.2 | 14.5 | 14.4 | 14.7 | 59.2 | 62.2 | 61.5 | 60.9 | 21.4 | 20.9 | 20.5 | 20.9 | −4.72% |
| *Coh* | 38.4 | 50.3 | 44.0 | 44.2 | 63.7 | 61.4 | 53.6 | 59.6 | 40.7 | 44.5 | 38.0 | 41.0 | 15.3 | 14.7 | 14.5 | 14.8 | 60.2 | 61.8 | 61.1 | 61.0 | 21.7 | 21.8 | 22.4 | 22.0 | −3.66% |
| *Und* | 39.5 | 48.8 | 49.3 | 45.9 | 49.1 | 57.7 | 58.5 | 55.1 | 43.8 | 45.1 | 42.2 | 43.7 | 14.9 | 14.4 | 14.7 | 14.7 | 58.3 | 61.8 | 62.2 | 60.8 | 21.1 | 20.3 | 20.9 | 20.8 | −4.44% |
| *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR** | 53.5 | 51.4 | 47.2 | 50.7 | 67.3 | 65.9 | 65.4 | 66.2 | 45.9 | 45.0 | 45.4 | 45.4 | 16.4 | 13.6 | 14.2 | 14.7 | 62.1 | 61.6 | 61.7 | 61.8 | 21.6 | 21.3 | 22.5 | 21.8 | +2.61% |
| *LESS* | 35.9 | 41.9 | 54.2 | 44.0 | 56.9 | 21.1 | 72.7 | 50.2 | 42.9 | 41.1 | 45.2 | 43.0 | 13.9 | 14.4 | 12.8 | 13.7 | 65.0 | 64.2 | 68.0 | 65.7 | 18.0 | 19.4 | 20.9 | 19.5 | −7.35% |
| *Applying LangGPS (to random selection and three best well-performing baselines marked with *)* | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 48.8 | 52.0 | 54.9 | 51.9 (+0.8) | 59.1 | 72.2 | 73.2 | 68.2 (+0.5) | 48.8 | 46.2 | 44.2 | 46.4 (+0.03) | 15.4 | 14.8 | 14.5 | 14.9 (-0.1) | 64.7 | 65.0 | 64.9 | **64.9** (+0.5) | 22.0 | 21.7 | 22.4 | 22.0 (-0.3) | +4.33% |
| *KMC* | 50.3 | 53.7 | 53.6 | 52.5 (+1.7) | 78.8 | 70.1 | 72.2 | **73.7** (+3.1) | 43.9 | 46.5 | 45.2 | 45.2 (+0.4) | 13.8 | 14.4 | 13.9 | 14.0 (+0.4) | 62.7 | 62.3 | 63.4 | 62.8 (+0.9) | 21.6 | 23.0 | 23.0 | **22.5** (+0.6) | +4.32% |
| *MTLD* | 51.8 | 52.7 | 53.7 | **52.7** (+0.9) | 69.8 | 70.5 | 68.4 | 69.6 (-0.3) | 47.2 | 47.6 | 43.3 | 46.1 (+0.7) | 15.2 | 15.0 | 14.9 | 15.0 (+0.3) | 65.5 | 64.3 | 63.9 | 64.6 (+0.7) | 21.1 | 21.8 | 22.4 | 21.8 (-0.4) | **+4.64%** |
| *DSIR* | 49.8 | 54.1 | 50.6 | 51.5 (+0.8) | 62.1 | 67.2 | 69.8 | 66.4 (+0.2) | 48.9 | 50.2 | 46.9 | **48.7** (+3.3) | 15.9 | 14.8 | 15.0 | **15.2** (+0.5) | 64.2 | 63.7 | 64.1 | 64.0 (+2.2) | 21.0 | 22.2 | 22.5 | 21.9 (+0.1) | +4.62% |

Table 1: Main results averaged over all languages for each dataset. **Avg.** denotes the average performance across the 1%, 3%, and 5% training settings. For each model, *LangGPS* is applied on top of both random selection and the three best well-performing baselines (with highest Δ values) under that model. Blue cell indicates better performance than the vanilla baseline under the same training setting, while Gray cell indicate a performance drop. Underline numbers indicate the best performance among vanilla baselines, **Bold** numbers indicate the best performance achieved after applying *LangGPS*. Δ represents the average relative gain or decline (in percentage) of each method across all datasets compared to training on the full dataset.

$$\Delta(\text{Method}) = \frac{1}{|\text{AllDatasets}|}\sum_{d\in\text{AllDatasets}} \frac{\text{Result(Method},d) - \text{Result(Full},d)}{\text{Result(Full},d)}$$

For each dataset, we evaluate on 10 languages, covering a total of 22 languages. *Accuracy* metric is used for *XNLI*, *XStoryCloze*, *MMMLU*, *MKQA* and *XQuAD* datasets. And for *XLSum* dataset, *ROUGE-L* scores are reported. We use greedy decoding with a max of 40 new tokens for each model. More details are provided in Appendix A.3,A.4.

**Baselines.** Detailed implementations are in Appendix A.5.

- **Random** *(Rand)*: samples are randomly selected from the full training set. Results are averaged over three different random seeds to ensure robustness.
- **KMeans Clustering** *(KMC)* performs kmeans in the embedding space of model `all-MiniLM-L6-v2` and selects the closest sample to each cluster centroid.
- **Lexical Diversity** *(MTLD)* (McCarthy and Jarvis 2010): selects samples with higher lexical diversity.
- **Naturalness** *(Nat)* (Zhong et al. 2022) selects samples that better resemble natural human-written text.
- **Coherence** *(Coh)* (Zhong et al. 2022) selects samples where the response serves as a coherent continuation of the paired question.
- **Understandability** *(Und)* (Zhong et al. 2022) selects samples that are more understandable.
- **Importance Resampling** *(DSIR)* (Xie et al. 2023) selects samples most relevant to the target set by estimating their importance weights.
- **Gradient Similarity** *(LESS)* (Xia et al. 2024): selects samples most relevant to the target set by low-rank gradient similarity search.

The target set for the target-dependent baselines *DSIR* and *LESS* is constructed by sampling 3 examples from each language subset across all 6 benchmarks.

**Detailed Implementations of *LangGPS*.** As a pre-selection strategy, *LangGPS* first selects the top $\rho\%$ of samples with the highest language separability scores from each language cluster in the original multilingual training set (with $\rho = 20$ in our main experiments). These selected samples are then passed to existing data selection methods for further fine-grained filtering. In our main experiments, we apply *LangGPS* on top of both *Random* selection and the three best-performing baselines for each model. We consider data selection settings of 1%, 3%, and 5%—i.e., the final training set contains 1%, 3%, or 5% of the full data—and report the average performance across these three settings.

### 4.2 Main results

Table 1 reports average results across languages for each dataset, and Figure 2 shows the average relative gains or declines on high- and low-resource languages. The detailed results for each language can be found in Appendix A.6.

**(1) Existing Methods Are Not Universally Effective.** Existing feature-based and target-dependent data selection methods fail to deliver consistent improvements in joint multilingual instruction tuning. For example, the trending method *LESS* achieves the best performance on *LLaMA-3.1-8B* ($\Delta = +4.90\%$), yet ranks the worst on *Qwen2.5-7B* ($\Delta = -7.35\%$). Moreover, *Nat*, *Coh*, and *Und*—the three



Figure 2: Average relative gains or declines (in percentage) when applying *LangGPS* on top of *Rand* and the three strongest-performing baselines, shown separately for high-resource and low-resource languages.

text quality-based selection methods—perform significantly worse on understanding tasks compared to generation tasks. And on *Qwen2.5-7B*, we observe that all baselines underperform even the *Random* selection strategy. These observations highlight the challenge of designing universally effective selection strategies in multilingual settings.

**(2) *LangGPS* Yields a Performance Boost, Especially on Understanding Tasks and Low-Resource Languages.** Overall, applying *LangGPS* on top of existing selection methods achieves consistent performance improvements, as evidenced by the increase in $\Delta$ values. And the gains are more pronounced on understanding tasks. Besides, we also observe that for *LLaMA-3.1-8B*, which exhibits weaker inherent multilingual capabilities, both the baselines and *LangGPS* yield more substantial improvements compared to *Qwen2.5-7B*, whose multilingual competence is already stronger. In Figure 2, we present the average performance gains or declines when applying *LangGPS*, separately for high-resource and low-resource languages[1]. Overall, *LangGPS* brings positive improvements on both language groups and the gains are more pronounced on low-resource languages, indicating that *LangGPS* can help mitigate the performance gap between high- and low-resource languages.

## 5 Further Analysis

### 5.1 Analysis of the Pre-selection Ratio $\rho$

We vary the pre-selection ratio $\rho$ from 10% to 100% and report the performance on *MMMLU*, *XLSum* in Figure 3

---

[1]We follow the taxonomy in https://microsoft.github.io/linguisticdiversity/assets/lang2tax.txt, where languages rated 4 or 5 are considered relatively high-resource, and others low-resource. And since the 10 languages selected from *MKQA* are not typical low-resource, *MKQA* dataset is excluded from the figure.

Figure 3: Effect of the pre-selection ratio $\rho$ on performance (*MMMLU* and *XLSum*, 5% setting, *LLaMA-3.1-8B*). We vary $\rho$ from 10% to 100% and report results using both *Rand* and *LESS* as downstream selectors.

(5% setting, *LLaMA-3.1-8B*) using both *Rand* and *LESS* as downstream selectors. Note that when $\rho = 100\%$, *LangGPS+baseline* becomes equivalent to the vanilla baseline.

We observe that *LangGPS* is not highly sensitive to the choice of $\rho$, and generally performs well when $\rho$ is in the range of 20% to 70%. A very small $\rho$ (e.g., 10%) tends to over-focus on highly separable samples, le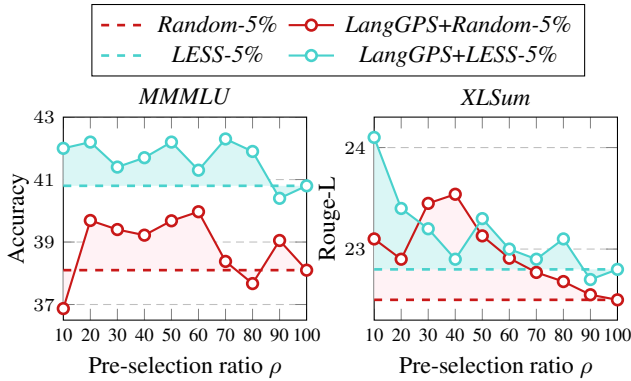ading to reduced diversity (as shown in Figure 1 (b)), and potential performance drops. On the other hand, as $\rho$ approaches 100%, the pre-selection effect of *LangGPS* gradually converges toward those of the vanilla baselines. Moreover, the choice of $\rho$ also presents a trade-off between data diversity and selection efficiency. A smaller $\rho$ means that subsequent selection methods—especially those with higher computational cost (e.g., *LESS*)—operate on a much smaller candidate pool, reducing overhead at the cost of potentially reduced diversity.

**Note that computational cost is analyzed in Appendix B.**

## 5.2 Analysis of Linguistic Boundary Representation

To provide a more intuitive understanding of *LangGPS*, we employ t-SNE (Van der Maaten and Hinton 2008) to visualize the representations of 200 sentences sampled from *XNLI* in parallel across English, Chinese and Arabic.

As shown in Figure 5 (1) (2) (3), multilingual SFT leads to clearer linguistic boundaries and more compact clusters in the model's representation space, reflected by increased average silhouette scores. Moreover, comparing Figure 5 (2) vs. (4) and (3) vs. (5), we observe that applying *LangGPS* further amplifies this effect: models trained with highly separable samples exhibit even higher silhouette scores, suggesting that *LangGPS* helps the model better form and maintain clear language boundaries.

## 5.3 The Role of Low-Separability Samples

Our experiments show that prioritizing highly separable samples during training improves multilingual performance. But *are low-separability samples entirely unhelpful?*



Figure 4: Average translation performance of *LLaMA-3.1-8B* on *FLORES+* ($X \to En$ and $En \to X$, where $X =$ De, Fr, Zh) under varying training sizes, comparing models trained on top separability-scoring, bottom separability-scoring, and randomly selected samples.

Low-separability samples often entangle with other language clusters in representation space, lacking clear language-specific features—but they may instead serve as implicit bridges for cross-lingual alignment and information sharing. To investigate this, we revisit the setup in Figure 1 (a) and Appendix A.1, training models with (1) top separability-scoring, (2) bottom separability-scoring, and (3) randomly selected samples (with varying sizes), and evaluate them on the *FLORES+* translation task to assess the role of low-separability data in cross-lingual scenarios.

The results in Figure 4, together with Figure 1(a), reveal several observations in common: **(1)** The three settings differ notably when training with limited samples, but their performances converge as training data increases. **(2)** Low-separability samples lead to a severe cold-start problem—when trained on just 50 examples, models perform significantly worse than those trained on high-separability or random data, likely due to the lack of clear language-specific signals in a small set of entangled samples.

Furthermore, we observe that the trends differ significantly between the "$X \to En$" and "$En \to X$" translation directions. **(1)** In "$X \to En$", none of the training settings yields improvement in the stable phase (red shaded area)—suggesting that vanilla multilingual SFT offers limited benefit for $X \to En$ translation. Nevertheless, low-separability samples result in the least degradation and occasionally yield marginal gains (e.g., at 400, 500, or 3000 samples). **(2)** In contrast, the overall trend of "En $\to$ X" direction

| Silhouette Score = **0.6774** | Silhouette Score = **0.7762** | Silhouette Score = **0.7665** | Silhouette Score = **0.8081** | Silhouette Score = **0.7837** |



| (1) Before SFT | (2) SFT with *"Random-5%"* | (3) SFT with *"LESS-5%"* | (4) SFT with *"LangGPS + Random-5%"* | (5) SFT with *"LangGPS + LESS-5%"* |

Figure 5: t-SNE visualizations of output representations from *LLaMA-3.1-8B* under different training settings: before SFT, after SFT with *Rand-5%*, *LESS-5%*, *LangGPS+Rand-5%*, and *LangGPS+LESS-5%*. Average silhouette scores are also reported.

| Datasets | XNLI | XStory Cloze | MMMLU | MKQA | XQuAD | XLSum |
|---|---|---|---|---|---|---|
| Strategy | | | *Model: LLaMA-3.1-8B* | | | |
| *Random* | 33.9 | 56.8 | **36.8** | 11.1 | 59.8 | 22.1 |
| *Ascending* | **37.9** | 67.5 | 32.5 | 11.0 | 58.7 | 21.3 |
| *Descending* | 33.1 | 63.4 | 33.2 | 10.8 | 61.0 | **22.2** |
| *Balanced* | 36.8 | **68.8** | 34.4 | **12.1** | 61.0 | 22.0 |
| Strategy | | | *Model: Qwen2.5-7B* | | | |
| *Random* | 45.3 | 74.3 | 48.6 | 11.8 | 61.4 | 21.5 |
| *Ascending* | 36.4 | **76.1** | 44.6 | 11.5 | 61.8 | 21.5 |
| *Descending* | 44.5 | 61.6 | 45.1 | **12.2** | 61.7 | 21.9 |
| *Balanced* | **50.8** | 75.9 | **49.7** | 11.9 | **62.7** | **22.3** |

Table 2: Results of multilingual curriculum learning. We compare the three curriculum strategies—*Ascending*, *Descending*, and *Balanced*—with *Randomly Shuffled*.

shows clear gains from multilingual SFT. And within the 200-2000 sample dip (the cyan-shaded area), model trained with low-separability samples experiences a milder decline compared to those trained with high-separability samples.

**Summary.** With limited training data, low-separability samples—due to their entangled representations across languages—may contribute more to cross-lingual alignment rather than supporting language-specific modeling, whereas high-separability samples are more effective for learning language-specific features and enabling a better warm start.
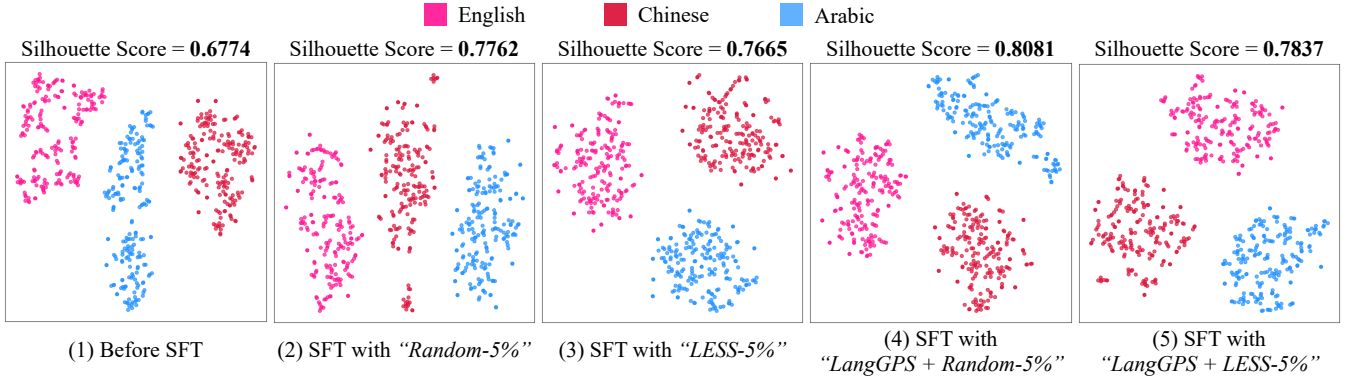
### 5.4 Language Separability Guided Multilingual Curriculum Learning

Curriculum learning shares an underlying intuition with data selection: both aim to improve model training by leveraging differences in sample utility. While data selection focuses on identifying and retaining a subset of training samples deemed beneficial, curriculum learning instead retains the full dataset but modulates the order in which the samples are presented. In essence, both strategies prioritize how the model interacts with training data—either by selecting what to learn from, or by deciding when to learn from it.

In this section, we explore whether language separability can serve as a guiding signal for multilingual curriculum learning. We design three curriculum strategies based on language separability: (1) *Ascending*: training progresses from low to high separability samples; (2) *Descending*: from high to low separability; (3) *Balanced*: training data are interleaved across different separability ranges to maintain an approximately uniform separability distribution throughout training (detailed implementations and algorithm are presented in Appendix C).

The results in Table 2 show that naively presenting data in strictly ascending or descending separability order leads to unsatisfactory performance. Instead, the *Balanced* strategy, where samples of varying separability are mixed evenly throughout training, yields the most stable and generalizable improvements across different models and datasets. This finding suggests that maintaining a diversity of separability levels during training helps the model benefit from both highly structured (language-specific) and entangled (cross-lingual) samples, striking a better trade-off.

## 6 Conclusion

We present *LangGPS*, a lightweight and broadly compatible data pre-selection framework that leverages language separability as the guiding signal for multilingual data selection. By quantifying how well different languages are distinguished in the model's representation space, *LangGPS* prioritizes highly separable samples that help models form clearer linguistic boundaries. Extensive experiments show that applying *LangGPS* on top of existing selection methods improves their effectiveness and generalizability in multilingual training, especially for understanding tasks and low-resource languages. Further analysis highlights the complementary roles of high- and low-separability samples, and demonstrates that language separability also benefits multilingual curriculum learning. Our findings position language separability as a principled signal for organizing multilingual data, offering a new perspective for improving the effectiveness and generalization of multilingual LLMs.

## Acknowledgements

**Limitations of this work are discussed in Appendix D.**

## References

Artetxe, M.; Ruder, S.; and Yogatama, D. 2020. On the Cross-lingual Transferability of Monolingual Representations. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 4623–4637. Online: Association for Computational Linguistics.

Bukharin, A.; Li, S.; Wang, Z.; Yang, J.; Yin, B.; Li, X.; Zhang, C.; Zhao, T.; and Jiang, H. 2024. Data Diversity Matters for Robust Instruction Tuning. In Al-Onaizan, Y.; Bansal, M.; and Chen, Y.-N., eds., *Findings of the Association for Computational Linguistics: EMNLP 2024*, 3411–3425. Miami, Florida, USA: Association for Computational Linguistics.

Che, W.; Feng, Y.; Qin, L.; and Liu, T. 2020. N-LTP: An open-source neural language technology platform for Chinese. *arXiv preprint arXiv:2009.11616*.

Chen, P.; Ji, S.; Bogoychev, N.; Kutuzov, A.; Haddow, B.; and Heafield, K. 2024. Monolingual or Multilingual Instruction Tuning: Which Makes a Better Alpaca. In Graham, Y.; and Purver, M., eds., *Findings of the Association for Computational Linguistics: EACL 2024*, 1347–1356. St. Julian's, Malta: Association for Computational Linguistics.

Conneau, A.; Rinott, R.; Lample, G.; Williams, A.; Bowman, S. R.; Schwenk, H.; and Stoyanov, V. 2018. XNLI: Evaluating Cross-lingual Sentence Representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.

Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Yang, A.; Fan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Hasan, T.; Bhattacharjee, A.; Islam, M. S.; Mubasshir, K.; Li, Y.-F.; Kang, Y.-B.; Rahman, M. S.; and Shahriyar, R. 2021. XL-Sum: Large-Scale Multilingual Abstractive Summarization for 44 Languages. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 4693–4703. Online: Association for Computational Linguistics.

Hendrycks, D.; Burns, C.; Basart, S.; Zou, A.; Mazeika, M.; Song, D.; and Steinhardt, J. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.

Huo, W.; Feng, X.; Huang, Y.; Fu, C.; Li, B.; Ye, Y.; Zhang, Z.; Tu, D.; Tang, D.; Lu, Y.; et al. 2025. Enhancing Non-English Capabilities of English-Centric Large Language Models Through Deep Supervision Fine-Tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 24185–24193.

Li, H.; Koto, F.; Wu, M.; Aji, A. F.; and Baldwin, T. 2023. Bactrian-x: Multilingual replicable instruction-following models with low-rank adaptation. *arXiv preprint arXiv:2305.15011*.

Li, M.; Zhang, Y.; Li, Z.; Chen, J.; Chen, L.; Cheng, N.; Wang, J.; Zhou, T.; and Xiao, J. 2024. From Quantity to Quality: Boosting LLM Performance with Self-Guided Data Selection for Instruction Tuning. In Duh, K.; Gomez, H.; and Bethard, S., eds., *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, 7602–7635. Mexico City, Mexico: Association for Computational Linguistics.

Lin, X. V.; Mihaylov, T.; Artetxe, M.; Wang, T.; Chen, S.; Simig, D.; Ott, M.; Goyal, N.; Bhosale, S.; Du, J.; Pasunuru, R.; Shleifer, S.; Koura, P. S.; Chaudhary, V.; O'Horo, B.; Wang, J.; Zettlemoyer, L.; Kozareva, Z.; Diab, M.; Stoyanov, V.; and Li, X. 2022. Few-shot Learning with Multilingual Generative Language Models. In Goldberg, Y.; Kozareva, Z.; and Zhang, Y., eds., *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 9019–9052. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics.

Longpre, S.; Lu, Y.; and Daiber, J. 2021. MKQA: A Linguistically Diverse Benchmark for Multilingual Open Domain Question Answering. *Transactions of the Association for Computational Linguistics*, 9: 1389–1406.

Lu, K.; Yuan, H.; Yuan, Z.; Lin, R.; Lin, J.; Tan, C.; Zhou, C.; and Zhou, J. 2023. # instag: Instruction tagging for analyzing supervised fine-tuning of large language models. *arXiv preprint arXiv:2308.07074*.

McCarthy, P. M.; and Jarvis, S. 2010. MTLD, vocd-D, and HD-D: A validation study of sophisticated approaches to lexical diversity assessment. *Behavior research methods*, 42(2): 381–392.

Muennighoff, N.; Wang, T.; Sutawika, L.; Roberts, A.; Biderman, S.; Scao, T. L.; Bari, M. S.; Shen, S.; Yong, Z.-X.; Schoelkopf, H.; et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.

Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744.

Qin, L.; Chen, Q.; Zhou, Y.; Chen, Z.; Li, Y.; Liao, L.; Li, M.; Che, W.; and Yu, P. S. 2024. Multilingual large language model: A survey of resources, taxonomy and frontiers. *arXiv preprint arXiv:2404.04925*.

Ranaldi, L.; Pucci, G.; and Freitas, A. 2024. Empowering cross-lingual abilities of instruction-tuned large language

models by translation-following demonstrations. In Ku, L.-W.; Martins, A.; and Srikumar, V., eds., *Findings of the Association for Computational Linguistics: ACL 2024*, 7961–7973. Bangkok, Thailand: Association for Computational Linguistics.

Rousseeuw, P. J. 1987. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20: 53–65.

Shaham, U.; Herzig, J.; Aharoni, R.; Szpektor, I.; Tsarfaty, R.; and Eyal, M. 2024. Multilingual instruction tuning with just a pinch of multilinguality. *arXiv preprint arXiv:2401.01854*.

Singh, S.; Vargus, F.; Dsouza, D.; Karlsson, B. F.; Mahendiran, A.; Ko, W.-Y.; Shandilya, H.; Patel, J.; Mataciunas, D.; OMahony, L.; Zhang, M.; Hettiarachchi, R.; Wilson, J.; Machado, M.; Moura, L. S.; Krzemiński, D.; Fadaei, H.; Ergün, I.; Okoh, I.; Alaagib, A.; Mudannayake, O.; Alyafeai, Z.; Chien, V. M.; Ruder, S.; Guthikonda, S.; Alghamdi, E. A.; Gehrmann, S.; Muennighoff, N.; Bartolo, M.; Kreutzer, J.; Üstün, A.; Fadaee, M.; and Hooker, S. 2024. Aya Dataset: An Open-Access Collection for Multilingual Instruction Tuning. arXiv:2402.06619.

Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Van der Maaten, L.; and Hinton, G. 2008. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11).

Xia, M.; Malladi, S.; Gururangan, S.; Arora, S.; and Chen, D. 2024. LESS: Selecting Influential Data for Targeted Instruction Tuning. In *International Conference on Machine Learning (ICML)*.

Xie, S. M.; Santurkar, S.; Ma, T.; and Liang, P. S. 2023. Data selection for language models via importance resampling. *Advances in Neural Information Processing Systems*, 36: 34201–34227.

Xu, C.; Sun, Q.; Zheng, K.; Geng, X.; Zhao, P.; Feng, J.; Tao, C.; Lin, Q.; and Jiang, D. 2024. WizardLM: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*.

Yang, A.; Yang, B.; Zhang, B.; Hui, B.; Zheng, B.; Yu, B.; Li, C.; Liu, D.; Huang, F.; Wei, H.; Lin, H.; Yang, J.; Tu, J.; Zhang, J.; Yang, J.; Yang, J.; Zhou, J.; Lin, J.; Dang, K.; Lu, K.; Bao, K.; Yang, K.; Yu, L.; Li, M.; Xue, M.; Zhang, P.; Zhu, Q.; Men, R.; Lin, R.; Li, T.; Xia, T.; Ren, X.; Ren, X.; Fan, Y.; Su, Y.; Zhang, Y.; Wan, Y.; Liu, Y.; Cui, Z.; Zhang, Z.; and Qiu, Z. 2024. Qwen2.5 Technical Report. *arXiv preprint arXiv:2412.15115*.

Yang, X.; Nie, S.; Liu, L.; Gururangan, S.; Karn, U.; Hou, R.; Khabsa, M.; and Mao, Y. 2025. Diversity-driven data selection for language model tuning through sparse autoencoder. *arXiv preprint arXiv:2502.14050*.

Ye, J.; Tao, X.; and Kong, L. 2023. Language versatilists vs. specialists: An empirical revisiting on multilingual transfer ability. *arXiv preprint arXiv:2306.06688*.

Ye, Y.; Feng, X.; Feng, X.; Ma, W.; Qin, L.; Xu, D.; Yang, Q.; Liu, H.; and Qin, B. 2024a. GlobeSumm: A Challenging Benchmark Towards Unifying Multi-lingual, Cross-lingual and Multi-document News Summarization. In Al-Onaizan, Y.; Bansal, M.; and Chen, Y.-N., eds., *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 10803–10821. Miami, Florida, USA: Association for Computational Linguistics.

Ye, Y.; Feng, X.; Feng, X.; Qin, L.; Huang, Y.; Huang, L.; Ma, W.; Hong, Q.; Zhang, Z.; Lu, Y.; et al. 2024b. Exploring Cross-lingual Latent Transplantation: Mutual Opportunities and Open Challenges. *arXiv preprint arXiv:2412.12686*.

Ye, Y.; Feng, X.; Yuan, Z.; Feng, X.; Qin, L.; Huang, L.; Ma, W.; Huang, Y.; Zhang, Z.; Lu, Y.; Yan, X.; Tang, D.; Tu, D.; and Qin, B. 2025. CC-Tuning: A Cross-Lingual Connection Mechanism for Improving Joint Multilingual Supervised Fine-Tuning. arXiv:2506.00875.

Zhang, X.; Li, S.; Hauer, B.; Shi, N.; and Kondrak, G. 2023. Don't Trust ChatGPT when your Question is not in English: A Study of Multilingual Abilities and Types of LLMs. In Bouamor, H.; Pino, J.; and Bali, K., eds., *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 7915–7927. Singapore: Association for Computational Linguistics.

Zhang, Y.; Wang, Y.; Liu, Z.; Wang, S.; Wang, X.; Li, P.; Sun, M.; and Liu, Y. 2024. Enhancing Multilingual Capabilities of Large Language Models through Self-Distillation from Resource-Rich Languages. *arXiv preprint arXiv:2402.12204*.

Zhao, H.; Andriushchenko, M.; Croce, F.; and Flammarion, N. 2024. Long is more for alignment: A simple but tough-to-beat baseline for instruction fine-tuning. *arXiv preprint arXiv:2402.04833*.

Zhao, Y.; Du, L.; Ding, X.; Ouyang, Y.; Wang, H.; Xiong, K.; Gao, J.; Sun, Z.; Xu, D.; Qing, Y.; et al. 2025. Beyond similarity: A gradient-based graph method for instruction tuning data selection. *arXiv preprint arXiv:2502.11062*.

Zheng, Y.; Zhang, R.; Zhang, J.; Ye, Y.; Luo, Z.; Feng, Z.; and Ma, Y. 2024. LlamaFactory: Unified Efficient Fine-Tuning of 100+ Language Models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*. Bangkok, Thailand: Association for Computational Linguistics.

Zhong, M.; Liu, Y.; Yin, D.; Mao, Y.; Jiao, Y.; Liu, P.; Zhu, C.; Ji, H.; and Han, J. 2022. Towards a Unified Multi-Dimensional Evaluator for Text Generation. In Goldberg, Y.; Kozareva, Z.; and Zhang, Y., eds., *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2023–2038. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics.

Zhou, C.; Liu, P.; Xu, P.; Iyer, S.; Sun, J.; Mao, Y.; Ma, X.; Efrat, A.; Yu, P.; Yu, L.; et al. 2023. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36: 55006–55021.

Zhu, W.; Lv, Y.; Dong, Q.; Yuan, F.; Xu, J.; Huang, S.; Kong, L.; Chen, J.; and Li, L. 2023. Extrapolating large language

models to non-english by aligning languages. *arXiv preprint arXiv:2308.04948.*

# A  Experiment Details

## A.1  Preliminary Experiments

We conduct preliminary experiments by fine-tuning *LLaMA-3.1-8B* using three types of training samples: (1) **High-separability** samples with the highest separability scores, (2) **Low-separability** samples with the lowest separability scores, and (3) **Random** selected samples from the full dataset. **To ensure fair comparison**, we select samples in accordance with the original language distribution of the dataset. For the random setting, we perform three independent samplings and report the averaged results. We vary the training data size from 0 to 5,000 to observe the impact of separability under different data scales. The full training corpus and *MMMLU* dataset used in preliminary experiments follows the settings in Appendix A.2,A.3.

Besides, we also present t-SNE visualizations of samples selected by different strategies under the same language distribution, using the last-token hidden states from *LLaMA-3.1-8B* as representations.

## A.2  Training Corpus

We totally select 97,696 multilingual instruction pairs from *aya dataset* (Singh et al. 2024) as our full training corpus and the training corpus covers 31 languages, ensuring extensive multilingual coverage. Detailed statistics are as follows:

```
1  {
2      "Basque": 939,
3      "Bengali": 1534,
4      "Egyptian Arabic": 529,
5      "English": 3944,
6      "French": 1422,
7      "German": 241,
8      "Greek": 623,
9      "Haitian": 106,
10     "Hindi": 1153,
11     "Indonesian": 786,
12     "Italian": 738,
13     "Japanese": 6259,
14     "Korean": 361,
15     "Moroccan Arabic": 8090,
16     "Najdi Arabic": 136,
17     "Portuguese": 8997,
18     "Russian": 423,
19     "Simplified Chinese": 3038,
20     "South Levantine Arabic": 81,
21     "Spanish": 3854,
22     "Standard Arabic": 4995,
23     "Swahili": 366,
24     "Ta'izzi-Adeni Arabic": 129,
25     "Tamil": 14133,
26     "Telugu": 8439,
27     "Thai": 724,
28     "Turkish": 4046,
29     "Ukrainian": 522,
30     "Urdu": 654,
31     "Vietnamese": 8676,
32     "Yoruba": 11758
33 }
```

## A.3  Evaluation Datasets

We conduct experiments on 6 benchmarks, which can be categorized into:

- **Multilingual Understanding:** (1) *XNLI* (Conneau et al. 2018), a multilingual natural language inference (NLI) dataset, (2) *XStoryCloze* (Lin et al. 2022), a multilingual commonsense reasoning dataset for evaluating story understanding and (3) *MMMLU*, the multilingual version of *MMLU* (Hendrycks et al. 2020), designed to evaluate models' general knowledge.
- **Multilingual Generation:** (1) *MKQA* (Longpre, Lu, and Daiber 2021), an open-domain multilingual question answering evaluation dataset, (2) *XQuAD* (Artetxe, Ruder, and Yogatama 2020), a question answering dataset and (3) *XLSum* (Hasan et al. 2021), a multilingual abstractive summarization benchmark comprising professionally annotated article-summary pairs.

For each dataset, we conduct experiments on 10 language, covering a total of 22 languages. For *XNLI*, *XStoryCloze*, *MMMLU*, *MKQA* and *XQuAD* datasets, *Accuracy* metric is used for evaluation. And for *XLSum* dataset, *ROUGE-L* scores are reported. We use greedy decoding with a max of 40 new tokens for each model.

*XNLI*, *XStoryCloze*, and *MMMLU* all belong to the multiple-choice category. For these datasets, a model's response is considered correct only if it contains the correct option and excludes all other options. For the short QA generative dataset *MKQA* and *XQuAD*, a model's answer is deemed correct if the gold answer appears in the model's response.

---

**Involved Languages (10 languages each dataset)**
*XNLI:* en, ar, el, es, fr, hi, sw, tr, vi, zh
*XStoryCloze:* en, ar, es, eu, hi, id, ru, sw, te, zh
*MMMLU:* en, ar, bn, de, es, fr, hi, sw, yo, zh
*MKQA:* en, ar, de, ja, ko, pt, ru, tr, vi, zh
*XQuAD:* en, ar, de, el, es, hi, ru, th, tr, zh
*XLSum:* en, es, fr, ko, pt, sw, tr, uk, vi, zh

---

**Sample Size**
*XNLI:* $1000 \times 10 = 10000$ (parallel)
*XStoryCloze:* $1511 \times 10 = 15110$ (parallel)
*MMMLU:* $1000 \times 10 = 10000$ (parallel)
*MKQA:* $1000 \times 10 = 10000$ (parallel)
*XQuAD:* $1190 \times 10 = 11900$ (parallel)
*XLSum:* $100 \times 10 = 1000$ (non-parallel)

---

## A.4  Prompts

The prompts we used for each datasets are listed in Table 4.

## A.5  Baselines Settings

The implementations of the baseline methods involved.

**Baselines.**

- **Random** *(Rand)*: samples are randomly selected from the full training set. Results are averaged over three different random seeds to ensure robustness.

| Composition | Getting Representations | | Clustering | | Data Selection | |
|---|---|---|---|---|---|---|
| | **Complexity** | **Actual** | **Complexity** | **Actual** | **Complexity** | **Actual** |
| *KMC* | $\mathcal{O}(|\mathcal{D}|)$ | 10 Mins | $\mathcal{O}(|\mathcal{D}| \cdot K)$ | 1 Min | $\mathcal{O}(|\mathcal{D}|)$ | 1 Min |

| Composition | MTLD Scoring | | Data Selection | |
|---|---|---|---|---|
| | **Complexity** | **Actual** | **Complexity** | **Actual** |
| *MTLD* | $\mathcal{O}(|\mathcal{D}|)$ | 10 Mins | $\mathcal{O}(|\mathcal{D}| \cdot \log|\mathcal{D}|)$ | 1 Min |

| Composition | UniEval Scoring | | Data Selection | |
|---|---|---|---|---|
| | **Complexity** | **Actual** | **Complexity** | **Actual** |
| *Nat/Coh/Und* | $\mathcal{O}(|\mathcal{D}|)$ | 2 Hours | $\mathcal{O}(|\mathcal{D}| \cdot \log|\mathcal{D}|)$ | 1 Min |

| Composition | Importance Weights Computation & Data Selection | |
|---|---|---|
| | **Complexity** | **Actual** |
| *DSIR* | - | 3 Mins |

| Composition | Warmup LoRA Training | | Gradient Features Computation | | Data Selection | |
|---|---|---|---|---|---|---|
| | **Complexity** | **Actual** | **Complexity** | **Actual** | **Complexity** | **Actual** |
| *LESS* | $\mathcal{O}(K)$ | 2 Hours | $\mathcal{O}(\mathcal{D})$ | 6 Hours | $\mathcal{O}(|\mathcal{D}| \cdot |\mathcal{D}_{tgt}|)$ | 1 Min |

| Composition | Getting Representations | | Silhouette Scoring | | Data Selection | |
|---|---|---|---|---|---|---|
| | **Complexity** | **Actual** | **Complexity** | **Actual** | **Complexity** | **Actual** |
| *LangGPS* | $\mathcal{O}(|\mathcal{D}|)$ | 2 Hours (fp16) | $\mathcal{O}(|\mathcal{D}|^2)$ | 1 Min | $\mathcal{O}(|\mathcal{D}| \cdot \log|\mathcal{D}|)$ | 1 Min |

Table 3: The computational cost of *LangGPS* and other baselines on full training corpus (measured as **single** A100 GPU hours, $|\mathcal{D}| = 97696$, $K = 5\% \times |\mathcal{D}| = 4885$ and $\mathcal{D}_{tgt}$ represents the target set for the target-dependent baselines). Getting Representations is the most costly step in *LangGPS*.

- **KMeans Clustering *(KMC)*** performs k-means in the embedding space of model `all-MiniLM-L6-v2`[2] and selects the closest sample to each cluster centroid. `all-MiniLM-L6-v2` maps sentences & paragraphs to a 384 dimensional dense vector space. To select $k$ representative samples, we perform k-means clustering in this embedding space and set the number of clusters in k-means to $k$. And then we identify the sample closest to each cluster centroid as the selected instance.

- **Lexical Diversity *(MTLD)*** (McCarthy and Jarvis 2010): selects samples with higher lexical diversity. We use LTP (Che et al. 2020) for word segmentation and the `lexical-diversity` Python package for computing MTLD scores, and select the top-scoring samples for training.

- **Naturalness *(Nat)*** (Zhong et al. 2022) selects samples that better resemble natural human-written text. We utilize `unieval-dialog` (Zhong et al. 2022) for computing naturalness scores, and select the top-scoring samples for training.

- **Coherence *(Coh)*** (Zhong et al. 2022) selects samples where the response serves as a coherent continuation of the paired question. We utilize `unieval-dialog` (Zhong et al. 2022) for computing coherence scores, and select the top-scoring samples for training.

- **Understandability *(Und)*** (Zhong et al. 2022) selects samples that are more understandable. We utilize `unieval-dialog` (Zhong et al. 2022) for computing understandability scores, and select the top-scoring samples for training.

- **Importance Resampling *(DSIR)*** (Xie et al. 2023) selects samples most relevant to the target set by estimating their importance weights.

- **Gradient Similarity *(LESS)*** (Xia et al. 2024): selects samples most relevant to the target set by low-rank gradient similarity search.

The target set for the target-dependent baselines ***DSIR*** and ***LESS*** is constructed by sampling 3 examples from each language subset across all 6 benchmarks.

---

[2]We use the English encoder `all-MiniLM-L6-v2` instead of a multilingual encoder because multilingual encoders are often trained to align semantically similar sentences across languages into a shared embedding space. While this design benefits cross-lingual retrieval or transfer learning tasks, it also has a notable side effect: they tend to suppress language-specific information. That is, samples from different languages but with similar semantics may be embedded closely together, making it difficult to distinguish languages based solely on their representations. In our context—where language-specific structure and diversity are important signals for data selection—such alignment behavior may be counterproductive. In contrast, monolingual models like `all-MiniLM-L6-v2` retain more surface-level and syntactic features, making them more suitable for preserving instructional variation and language-specific patterns.

## A.6 Detailed Experimental Results on Each Involved Language

In this section, we present the detailed results for each language involved in Table 5, 6, 7, 8, 9, 10.

# B Computational Cost Analysis

In this section, we report the actual time cost of *LangGPS* and other baselines in Table 3, measured as **single** A100 GPU hours on full training corpus.

We observe that the majority of *LangGPS*'s computational overhead comes from the *Getting Representation* step, where we pass the training samples into the model and get the representations of its last input token, while the time spent on *Silhouette Scoring* and *Data Selection* is negligible in practice. Notably, though *Getting Representation* costs hours, the process itself is quite simple and implementation-friendly—requiring only a single forward pass[3] over the training corpus without any additional computation.

Besides, as a pre-selection strategy, *LangGPS* can help reduce the computational overhead of subsequent selection methods—particularly those with higher costs, such as *LESS*—by narrowing down the candidate pool. Specifically, under a pre-selection ratio of $\rho = 20\%$, the total cost of *LangGPS+LESS* becomes:

$$
\begin{aligned}
\text{Cost}(\textit{LangGPS+LESS}) &= \text{Cost}(\textit{LangGPS}_{\text{total}}) \\
&+ \text{Cost}(\textit{LESS}_{\text{Warmup}}) \\
&+ \mathbf{20\%} \times \text{Cost}(\textit{LESS}_{\text{Gradient}}) \\
&\approx 5.2 \text{ Hours} < \text{Cost}(\textit{LESS}_{\text{total}}) = 8 \text{ Hours}
\end{aligned}
$$

# C Curriculum Learning Implementation

We first sort the entire training corpus by language separability scores and divide it into 10 equally sized buckets based on score percentiles: Top 0–10%, Top 10–20%, ..., Top 90–100%. Specifically, each bucket contains the samples within the corresponding range of separability scores from each language cluster in the original multilingual training set. Then, we apply the following curriculum learning strategies:

- **Ascending:** Train the model starting from the samples with the *lowest* separability scores and gradually proceed to those with the *highest*. (Training order: Top 90–100%, 80–90%, ..., 0–10%)
- **Descending:** Train the model starting from the samples with the *highest* separability scores and gradually proceed to those with the *lowest*. (Training order: Top 0–10%, 10–20%, ..., 90–100%)
- **Balanced:** Training data are interleaved across different separability ranges to maintain an approximately uniform separability distribution throughout training, ensuring no score range is over- or under-represented. The detailed algorithm is as follows:

---

[3]`model(**inputs, output_hidden_states=True)`

---

**Algorithm 1: Balanced Training Data Construction**

1: Initialize training set $\mathcal{D}_{\text{train}} \leftarrow [\ ]$
2: Sort all training samples by language separability scores.
3: Divide them into 10 equal-sized buckets: $\mathcal{B}_1$ (Top 0-10%), $\mathcal{B}_2$ (Top 10-20%), ..., $\mathcal{B}_{10}$ (Top 90-100%), where each $\mathcal{B}_i$ contains the corresponding score-range samples from each language cluster.
4: **while** any bucket $\mathcal{B}_i$ is non-empty **do**
5:     Initialize mini-batch $\mathcal{M} \leftarrow [\ ]$
6:     **for** $i = 1$ to 10 **do**
7:         **if** $\mathcal{B}_i$ is not empty **then**
8:             Sample one instance $x$ from $\mathcal{B}_i$
9:             Append $x$ to $\mathcal{M}$
10:         **end if**
11:     **end for**
12:     Shuffle $\mathcal{M}$ and append to $\mathcal{D}_{\text{train}}$
13: **end while**
14: **return** $\mathcal{D}_{\text{train}}$

---

# D Limitation

This work exhibits several limitations worth noting. Firstly, the quantification of language separability in *LangGPS* essentially reflects the model's own multilingual modeling capability. As multilingual capacities vary across models, *LangGPS* requires model-specific data selection, lacking the simplicity of a one-size-fits-all solution. Secondly, our experiments were conducted on *LLaMA-3.1-8B* and *Qwen2.5-7B*. While these models represent important milestones in open-source LLM development, the evaluation across more LLMs would improve the generalizability of our findings across the broader LLM ecosystem. Thirdly, in Section 5.3, we examine the impact of high- and low-separability samples on translation performance. It is worth noting that the training data (aya dataset) is not designed to improve translation quality. Hence, the goal of this section is not to show that low-separability samples enhance translation performance. Rather, we highlight that the milder decline of low-separability samples within a fluctuating performance region reveals their distinct functional role compared to high-separability samples.

| Prompt for *XNLI* (English version) |
| --- |

Premise: {premise}

Hypothesis: {hypothesis}

What do you think is the relationship between the premise and the hypothesis?

(1) Entail
(2) Neutral
(3) Contradict

If you have to choose one of these options, your answer would be: {response}

| Prompt for *XStoryCloze* (English version) |
| --- |

{story}

(1) {first possible continuation of the story}

(2) {second possible continuation of the story}

Which of the two options is more likely to be the ending of the given story.

Your answer: {response}

| Prompt for *MMMLU* (English version) |
| --- |

Please answer the following multiple choice question.

{question}

(A) {option 1}

(B) {option 2}

(C) {option 3}

(D) {option 4}

If you have to choose one of these options, your answer would be: {response}

| Prompt for *MKQA* (English version) |
| --- |

Please answer the following question.

{question}

Your answer: {response}

| Prompt for *XQuAD* (English version) |
| --- |

Please answer these questions only based on the given context.

Context: {context}

Question: {question}

Your answer: {response}

| Prompt for *XLSum* (English version) |
| --- |

Please summarize the following content into one sentence.

{content}

{response}

Table 4: The prompts used for *XNLI*, *XStoryCloze*, *MMMLU*, *MKQA*, *XQuAD* and *XLSum*.

**Model: LLaMA-3.1-8B**

**Dataset: XNLI**

| Method | 1% | | | | | | | | | | | 3% | | | | | | | | | | | 5% | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Langs | en | ar | el | es | fr | hi | sw | tr | vi | zh | **Avg.** | en | ar | el | es | fr | hi | sw | tr | vi | zh | **Avg.** | en | ar | el | es | fr | hi | sw | tr | vi | zh | **Avg.** |
| *Rand* | 35.4 | 10.0 | 23.1 | 34.6 | 30.0 | 10.6 | 29.3 | 31.4 | 36.2 | 10.5 | 25.1 | 35.1 | 30.0 | 30.7 | 35.7 | 36.6 | 30.7 | 30.0 | 31.2 | 39.7 | 30.4 | 33.0 | 40.6 | 31.2 | 33.2 | 40.9 | 39.6 | 30.7 | 32.4 | 35.1 | 37.3 | 32.5 | 35.4 |
| *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC\** | 28.8 | 29.7 | 26.5 | 30.7 | 17.2 | 18.5 | 29.5 | 23.7 | 36.2 | 1.7 | 24.3 | 37.8 | 32.7 | 31.0 | 37.3 | 35.2 | 31.4 | 34.9 | 36.4 | 41.9 | 36.4 | 35.5 | 19.9 | 31.7 | 30.4 | 44.8 | 40.9 | 33.7 | 30.6 | 36.7 | 37.2 | 34.8 | 34.1 |
| *MTLD* | 33.8 | 29.1 | 31.6 | 29.6 | 35.4 | 33.2 | 30.5 | 4.7 | 30.8 | 29.0 | 28.8 | 36.5 | 30.5 | 37.2 | 36.6 | 37.6 | 32.2 | 36.8 | 36.7 | 33.2 | 31.4 | 34.9 | 36.9 | 31.0 | 33.5 | 40.9 | 38.9 | 30.5 | 32.1 | 36.7 | 34.1 | 29.5 | 34.4 |
| *Nat* | 6.6 | 21.2 | 28.9 | 24.3 | 1.1 | 34.0 | 32.5 | 5.9 | 14.4 | 35.0 | 20.4 | 46.9 | 35.2 | 42.2 | 51.7 | 47.8 | 35.3 | 30.4 | 35.6 | 37.1 | 38.9 | 40.1 | 34.1 | 33.8 | 31.5 | 47.4 | 40.0 | 30.8 | 29.8 | 36.5 | 43.8 | 34.2 | 36.2 |
| *Coh\** | 29.8 | 29.5 | 29.7 | 21.7 | 29.3 | 29.5 | 28.8 | 24.7 | 30.4 | 29.6 | 28.3 | 39.8 | 32.1 | 34.1 | 40.3 | 45.7 | 29.8 | 35.5 | 39.0 | 36.0 | 38.7 | 37.1 | 32.6 | 31.5 | 33.1 | 42.3 | 44.3 | 31.9 | 30.9 | 36.4 | 40.4 | 39.7 | 36.3 |
| *Und* | 8.6 | 29.3 | 24.7 | 33.5 | 2.8 | 38.8 | 31.8 | 33.1 | 0.9 | 31.6 | 23.5 | 41.5 | 30.6 | 33.6 | 42.6 | 46.5 | 32.0 | 33.0 | 36.4 | 40.6 | 43.8 | 38.1 | 30.5 | 30.7 | 32.0 | 48.5 | 38.2 | 30.7 | 30.3 | 37.0 | 37.8 | 34.0 | 35.0 |
| *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR* | 37.5 | 10.4 | 30.7 | 41.7 | 36.4 | 14.1 | 31.0 | 35.7 | 35.5 | 22.0 | 29.5 | 10.3 | 10.0 | 31.8 | 35.3 | 32.1 | 20.2 | 33.3 | 35.8 | 41.1 | 26.6 | 27.7 | 37.5 | 33.8 | 32.3 | 40.8 | 41.9 | 37.9 | 34.5 | 36.8 | 33.8 | 39.7 | 36.9 |
| *LESS\** | 37.0 | 29.9 | 33.1 | 39.3 | 34.8 | 27.7 | 35.1 | 33.6 | 39.4 | 29.4 | 33.9 | 38.9 | 30.0 | 29.3 | 44.3 | 35.6 | 29.4 | 29.2 | 29.7 | 39.6 | 30.5 | 33.7 | 43.0 | 35.7 | 34.9 | 47.2 | 40.4 | 30.7 | 32.2 | 35.3 | 43.4 | 35.5 | 37.8 |
| *Applying LangGPS (to random selection and three best well-performing baselines marked with \*)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 35.3 | 10.3 | 20.3 | 34.9 | 33.9 | 15.5 | 30.8 | 33.4 | 39.1 | 10.9 | 26.4 | 36.6 | 32.3 | 34.2 | 38.6 | 36.8 | 32.3 | 32.1 | 35.6 | 36.5 | 32.4 | 34.7 | 37.6 | 36.9 | 36.6 | 41.4 | 38.1 | 31.4 | 33.3 | 35.9 | 36.8 | 33.9 | 36.2 |
| *KMC* | 29.4 | 28.6 | 24.7 | 30.4 | 31.7 | 19.3 | 29.2 | 8.0 | 31.2 | 28.8 | 26.1 | 41.2 | 37.6 | 34.6 | 44.2 | 42.1 | 37.7 | 33.3 | 35.0 | 35.3 | 38.8 | 38.0 | 39.1 | 36.1 | 35.5 | 42.6 | 43.5 | 35.8 | 31.7 | 36.2 | 38.1 | 37.0 | 37.6 |
| *Coh* | 35.4 | 14.3 | 35.8 | 45.0 | 39.3 | 27.0 | 31.8 | 36.0 | 40.4 | 31.5 | 33.7 | 16.3 | 30.2 | 35.3 | 38.0 | 32.9 | 36.1 | 29.9 | 32.8 | 40.3 | 36.3 | 32.8 | 37.2 | 37.4 | 36.4 | 41.9 | 39.4 | 36.6 | 38.6 | 37.1 | 36.6 | 39.8 | 38.1 |
| *LESS* | 37.2 | 33.3 | 32.8 | 41.2 | 37.7 | 32.4 | 34.9 | 30.5 | 41.9 | 34.3 | 35.6 | 40.7 | 34.5 | 41.2 | 47.0 | 40.0 | 34.5 | 33.5 | 29.0 | 37.4 | 40.5 | 37.8 | 37.1 | 35.6 | 37.1 | 40.3 | 38.5 | 36.6 | 37.5 | 36.8 | 38.5 | 38.5 | 37.7 |

**Model: Qwen2.5-7B**

**Dataset: XNLI**

| Method | 1% | | | | | | | | | | | 3% | | | | | | | | | | | 5% | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Langs | en | ar | el | es | fr | hi | sw | tr | vi | zh | **Avg.** | en | ar | el | es | fr | hi | sw | tr | vi | zh | **Avg.** | en | ar | el | es | fr | hi | sw | tr | vi | zh | **Avg.** |
| *Rand* | 81.3 | 31.2 | 26.0 | 74.5 | 31.7 | 32.1 | 24.4 | 49.0 | 54.8 | 50.8 | 45.6 | 87.1 | 43.6 | 39.4 | 78.7 | 63.0 | 39.6 | 19.8 | 48.7 | 55.0 | 54.8 | 53.0 | 87.1 | 53.0 | 41.8 | 78.6 | 62.7 | 41.4 | 25.9 | 46.5 | 52.5 | 57.3 | 54.7 |
| *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC\** | 84.3 | 50.0 | 36.8 | 75.4 | 51.6 | 32.9 | 29.7 | 47.4 | 52.4 | 50.9 | 51.1 | 79.9 | 54.8 | 44.3 | 75.2 | 58.7 | 39.8 | 29.7 | 42.9 | 17.0 | 56.0 | 49.8 | 82.5 | 46.8 | 34.9 | 71.5 | 57.4 | 40.4 | 27.0 | 43.8 | 53.4 | 55.4 | 51.3 |
| *MTLD\** | 88.6 | 48.8 | 44.9 | 83.0 | 45.0 | 35.3 | 4.6 | 46.7 | 50.2 | 53.6 | 50.1 | 88.3 | 50.5 | 40.1 | 78.4 | 61.0 | 37.7 | 28.6 | 50.3 | 51.1 | 53.4 | 53.9 | 84.0 | 49.4 | 35.8 | 75.5 | 55.4 | 36.1 | 26.6 | 47.7 | 50.7 | 53.4 | 51.5 |
| *Nat* | 74.1 | 50.3 | 34.0 | 77.0 | 31.7 | 28.7 | 25.5 | 8.7 | 2.8 | 46.6 | 37.9 | 86.2 | 42.6 | 45.3 | 66.1 | 42.7 | 35.6 | 31.9 | 4.9 | 24.0 | 44.9 | 42.4 | 86.4 | 50.6 | 48.9 | 70.0 | 46.4 | 38.8 | 3.1 | 45.7 | 44.2 | 56.4 | 49.1 |
| *Coh* | 79.6 | 49.2 | 7.7 | 73.5 | 52.9 | 23.1 | 5.9 | 11.0 | 33.7 | 47.8 | 38.4 | 86.5 | 52.0 | 46.3 | 78.8 | 54.1 | 36.7 | 17.9 | 24.1 | 51.2 | 55.1 | 50.3 | 72.4 | 48.4 | 48.6 | 74.2 | 43.2 | 34.5 | 18.4 | 20.1 | 29.4 | 50.4 | 44.0 |
| *Und* | 81.1 | 47.7 | 22.9 | 76.8 | 13.8 | 33.6 | 13.9 | 22.3 | 27.1 | 56.1 | 39.5 | 86.5 | 49.9 | 48.3 | 71.6 | 46.0 | 36.8 | 33.2 | 35.5 | 31.1 | 49.5 | 48.8 | 83.9 | 48.8 | 46.7 | 68.2 | 48.3 | 35.9 | 29.2 | 39.7 | 36.5 | 55.3 | 49.3 |
| *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR\** | 88.2 | 48.0 | 41.3 | 79.7 | 57.9 | 41.5 | 32.9 | 38.8 | 50.8 | 56.0 | 53.5 | 85.6 | 49.6 | 37.7 | 73.8 | 59.9 | 37.6 | 26.0 | 42.7 | 45.5 | 55.1 | 51.4 | 82.6 | 50.5 | 42.2 | 70.5 | 52.1 | 38.4 | 7.6 | 36.2 | 41.5 | 50.7 | 47.2 |
| *LESS* | 80.5 | 34.5 | 37.4 | 37.3 | 32.5 | 5.2 | 34.0 | 30.4 | 36.0 | 30.8 | 35.9 | 51.2 | 46.6 | 44.3 | 48.0 | 42.4 | 41.0 | 17.7 | 34.2 | 46.8 | 46.8 | 41.9 | 84.8 | 54.0 | 34.9 | 76.5 | 58.5 | 32.9 | 26.4 | 50.8 | 59.4 | 63.3 | 54.2 |
| *Applying LangGPS (to random selection and three best well-performing baselines marked with \*)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 75.0 | 45.1 | 31.8 | 72.8 | 52.6 | 38.7 | 31.3 | 47.5 | 49.8 | 43.2 | 48.8 | 86.6 | 42.5 | 44.9 | 77.3 | 55.1 | 39.4 | 30.6 | 45.6 | 52.2 | 45.6 | 52.0 | 87.5 | 52.2 | 43.9 | 79.1 | 60.8 | 38.9 | 26.4 | 49.7 | 52.7 | 57.5 | 54.9 |
| *KMC* | 82.0 | 48.9 | 41.4 | 67.3 | 46.8 | 37.9 | 30.2 | 48.7 | 43.7 | 55.9 | 50.3 | 89.1 | 55.4 | 48.5 | 82.9 | 62.7 | 34.6 | 29.4 | 37.3 | 43.6 | 53.6 | 53.7 | 86.5 | 45.9 | 39.0 | 76.8 | 57.6 | 40.1 | 28.9 | 48.3 | 53.9 | 59.3 | 53.6 |
| *MTLD* | 83.9 | 45.8 | 39.0 | 77.5 | 53.9 | 40.3 | 24.7 | 49.6 | 53.4 | 50.2 | 51.8 | 88.8 | 49.3 | 41.0 | 82.0 | 52.5 | 41.8 | 11.2 | 50.5 | 52.8 | 56.7 | 52.7 | 88.6 | 54.3 | 45.0 | 82.7 | 63.9 | 38.1 | 6.8 | 46.3 | 53.3 | 58.4 | 53.7 |
| *DSIR* | 86.9 | 22.9 | 29.3 | 81.1 | 46.7 | 47.0 | 24.7 | 51.3 | 50.9 | 56.7 | 49.8 | 87.5 | 54.7 | 45.8 | 78.3 | 52.5 | 39.4 | 30.3 | 43.5 | 50.5 | 58.6 | 54.1 | 89.4 | 50.0 | 44.8 | 81.7 | 54.0 | 41.9 | 15.7 | 47.1 | 53.4 | 27.9 | 50.6 |

Table 5: The detailed performance results of different language subset on *XNLI* dataset.

**Model: LLaMA-3.1-8B**

**Dataset: XStoryCloze**

| Method | 1% en | ar | es | eu | hi | id | ru | sw | te | zh | Avg. | 3% en | ar | es | eu | hi | id | ru | sw | te | zh | Avg. | 5% en | ar | es | eu | hi | id | ru | sw | te | zh | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Rand* | 51.1 | 48.6 | 54.1 | 26.9 | 55.9 | 46.3 | 49.7 | 28.0 | 19.6 | 56.7 | 43.7 | 67.3 | 60.3 | 63.7 | 49.2 | 68.1 | 64.2 | 74.1 | 44.7 | 45.1 | 77.6 | 61.4 | 80.9 | 43.9 | 57.1 | 52.1 | 68.0 | 74.2 | 73.8 | 44.4 | 38.7 | 77.6 | 61.1 |
| *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC\** | 41.6 | 18.9 | 49.8 | 24.2 | 64.3 | 55.1 | 25.6 | 19.4 | 59.5 | 83.0 | 44.1 | 67.8 | 61.6 | 73.6 | 52.1 | 69.3 | 72.0 | 69.3 | 37.8 | 55.1 | 78.0 | 63.7 | 88.6 | 76.5 | 55.5 | 53.1 | 63.5 | 79.0 | 53.3 | 49.8 | 54.1 | 82.7 | 65.6 |
| *MTLD* | 38.0 | 6.0 | 36.3 | 24.0 | 60.0 | 65.1 | 32.4 | 19.0 | 0.0 | 62.3 | 34.3 | 82.0 | 10.1 | 68.2 | 47.3 | 57.3 | 72.9 | 56.3 | 48.6 | 11.9 | 74.5 | 52.9 | 67.2 | 14.4 | 56.7 | 52.0 | 21.0 | 69.6 | 44.3 | 37.9 | 27.3 | 80.7 | 47.1 |
| *Nat* | 0.2 | 1.1 | 0.3 | 1.4 | 6.5 | 24.4 | 40.3 | 20.6 | 0.1 | 37.9 | 13.3 | 54.6 | 81.2 | 21.6 | 4.8 | 53.0 | 74.6 | 32.0 | 31.0 | 0.0 | 47.4 | 40.0 | 53.0 | 78.6 | 13.4 | 47.8 | 61.9 | 80.3 | 75.6 | 54.3 | 0.9 | 82.5 | 54.8 |
| *Coh\** | 56.4 | 49.9 | 66.0 | 19.8 | 54.3 | 39.4 | 45.2 | 31.0 | 0.9 | 57.7 | 42.1 | 80.6 | 66.6 | 70.0 | 38.0 | 51.6 | 47.5 | 30.2 | 31.3 | 51.3 | 78.1 | 54.5 | 84.0 | 68.2 | 47.6 | 33.0 | 63.3 | 49.2 | 27.8 | 36.8 | 53.1 | 51.8 | 51.5 |
| *Und* | 58.8 | 9.7 | 0.0 | 1.8 | 72.0 | 38.8 | 15.4 | 4.1 | 0.5 | 51.9 | 25.3 | 80.3 | 61.8 | 34.5 | 44.5 | 66.7 | 62.5 | 26.1 | 20.1 | 11.5 | 67.0 | 47.5 | 77.5 | 83.6 | 13.4 | 63.1 | 69.0 | 79.1 | 59.0 | 56.4 | 9.3 | 86.2 | 59.6 |
| *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR* | 53.9 | 42.7 | 67.0 | 17.7 | 49.6 | 52.3 | 25.3 | 31.5 | 1.0 | 25.5 | 36.7 | 20.5 | 44.5 | 58.7 | 18.7 | 36.4 | 59.9 | 40.6 | 42.8 | 7.9 | 67.2 | 39.7 | 49.8 | 63.3 | 60.7 | 35.5 | 37.1 | 65.5 | 58.7 | 43.5 | 24.2 | 56.5 | 49.5 |
| *LESS\** | 62.4 | 76.1 | 68.7 | 39.7 | 73.0 | 65.3 | 69.2 | 51.4 | 69.9 | 79.0 | 65.5 | 90.1 | 76.9 | 84.6 | 42.8 | 71.3 | 75.2 | 81.9 | 50.0 | 55.9 | 82.8 | 71.2 | 83.9 | 68.6 | 83.9 | 55.7 | 75.2 | 80.1 | 79.2 | 46.5 | 41.5 | 68.0 | 68.3 |
| *Applying LangGPS (to random selection and three best well-performing baselines marked with \*)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 53.2 | 42.8 | 53.4 | 30.8 | 56.6 | 64.2 | 60.0 | 41.2 | 17.9 | 66.7 | 48.7 | 89.2 | 71.5 | 84.1 | 58.3 | 68.9 | 77.9 | 81.6 | 54.9 | 56.5 | 82.1 | 72.5 | 87.0 | 71.0 | 85.0 | 56.8 | 69.9 | 78.0 | 84.1 | 54.7 | 54.8 | 83.8 | 72.5 |
| *KMC* | 59.6 | 39.1 | 54.3 | 14.4 | 56.8 | 67.6 | 42.6 | 22.4 | 0.3 | 66.8 | 42.4 | 80.5 | 61.8 | 60.7 | 45.5 | 74.5 | 82.5 | 84.8 | 31.8 | 59.7 | 81.1 | 66.3 | 86.0 | 68.6 | 81.2 | 52.3 | 69.6 | 73.6 | 71.5 | 53.5 | 54.7 | 81.4 | 69.2 |
| *Coh* | 32.4 | 57.2 | 82.2 | 34.9 | 61.6 | 67.0 | 31.2 | 24.6 | 2.8 | 77.0 | 47.1 | 22.2 | 75.4 | 86.8 | 40.0 | 31.1 | 72.5 | 24.2 | 37.3 | 4.6 | 77.8 | 47.2 | 27.1 | 76.6 | 86.0 | 41.1 | 68.4 | 69.4 | 64.2 | 41.8 | 63.5 | 70.2 | 60.8 |
| *LESS* | 87.2 | 70.6 | 81.9 | 45.4 | 62.7 | 69.8 | 73.7 | 29.8 | 58.5 | 81.2 | 66.1 | 89.8 | 82.5 | 87.0 | 57.4 | 79.9 | 80.3 | 87.0 | 61.2 | 61.0 | 82.3 | 76.8 | 90.7 | 78.4 | 86.0 | 53.1 | 78.6 | 81.5 | 85.1 | 52.3 | 61.0 | 73.0 | 74.0 |

**Model: Qwen2.5-7B**

**Dataset: XStoryCloze**

| Method | 1% en | ar | el | es | fr | hi | sw | tr | vi | zh | Avg. | 3% en | ar | el | es | fr | hi | sw | tr | vi | zh | Avg. | 5% en | ar | el | es | fr | hi | sw | tr | vi | zh | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Rand* | 60.6 | 50.6 | 41.6 | 57.8 | 61.4 | 55.4 | 54.8 | 31.4 | 55.2 | 73.5 | 54.2 | 90.2 | 88.6 | 90.8 | 58.1 | 76.7 | 85.3 | 90.2 | 50.4 | 50.9 | 78.5 | 76.0 | 83.5 | 85.9 | 90.1 | 52.2 | 74.3 | 83.9 | 89.6 | 43.4 | 50.4 | 74.2 | 72.7 |
| *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC\** | 90.5 | 84.1 | 91.7 | 58.2 | 78.7 | 85.8 | 86.6 | 27.1 | 56.5 | 83.7 | 74.3 | 78.1 | 81.1 | 57.4 | 50.6 | 66.0 | 76.5 | 61.9 | 54.2 | 39.7 | 78.6 | 64.4 | 84.5 | 86.0 | 90.5 | 56.3 | 74.0 | 76.4 | 74.3 | 52.5 | 58.9 | 76.4 | 73.0 |
| *MTLD\** | 83.2 | 58.0 | 80.7 | 52.9 | 65.5 | 80.1 | 82.4 | 31.4 | 54.3 | 53.3 | 64.2 | 84.3 | 85.7 | 83.8 | 57.0 | 74.7 | 86.8 | 87.2 | 49.1 | 50.6 | 76.3 | 73.6 | 89.0 | 83.5 | 80.3 | 56.9 | 74.1 | 86.7 | 70.9 | 50.1 | 47.3 | 80.8 | 72.0 |
| *Nat* | 85.6 | 27.0 | 88.2 | 29.1 | 80.5 | 71.8 | 25.9 | 47.2 | 44.8 | 66.6 | 56.7 | 46.7 | 83.4 | 86.4 | 6.6 | 80.9 | 75.8 | 27.5 | 27.7 | 0.5 | 77.8 | 51.3 | 63.3 | 56.1 | 91.9 | 43.5 | 60.4 | 86.4 | 44.5 | 17.1 | 39.0 | 85.8 | 58.8 |
| *Coh* | 81.1 | 82.5 | 81.3 | 45.7 | 64.2 | 80.7 | 68.2 | 13.6 | 44.6 | 74.9 | 63.7 | 73.1 | 78.1 | 82.7 | 39.7 | 72.9 | 79.7 | 55.8 | 30.5 | 46.3 | 55.5 | 61.4 | 63.5 | 70.5 | 72.3 | 46.9 | 68.9 | 76.5 | 43.5 | 10.2 | 47.2 | 36.6 | 53.6 |
| *Und* | 82.4 | 22.0 | 90.0 | 16.7 | 55.2 | 75.3 | 29.3 | 23.9 | 35.3 | 61.1 | 49.1 | 72.6 | 77.2 | 73.9 | 48.9 | 56.7 | 78.4 | 32.2 | 14.8 | 51.2 | 70.8 | 57.7 | 60.8 | 63.8 | 51.5 | 38.8 | 50.6 | 86.7 | 61.8 | 47.8 | 46.6 | 76.2 | 58.5 |
| *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR\** | 81.9 | 78.2 | 84.2 | 51.7 | 61.2 | 80.5 | 85.4 | 38.7 | 33.3 | 78.3 | 67.3 | 84.8 | 80.2 | 78.9 | 51.9 | 54.4 | 82.0 | 69.7 | 26.3 | 51.9 | 78.9 | 65.9 | 78.5 | 74.2 | 70.3 | 50.3 | 68.7 | 79.8 | 72.3 | 34.9 | 49.0 | 76.2 | 65.4 |
| *LESS* | 70.0 | 74.0 | 82.5 | 39.4 | 46.1 | 45.9 | 68.9 | 43.7 | 47.7 | 51.0 | 56.9 | 15.6 | 4.9 | 13.3 | 3.2 | 60.0 | 11.3 | 2.1 | 11.3 | 48.3 | 41.6 | 21.1 | 91.2 | 77.8 | 79.0 | 53.9 | 77.6 | 78.0 | 77.5 | 51.9 | 51.4 | 88.6 | 72.7 |
| *Applying LangGPS (to random selection and three best well-performing baselines marked with \*)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 59.1 | 55.0 | 48.8 | 49.3 | 62.7 | 70.7 | 67.0 | 42.9 | 55.8 | 80.0 | 59.1 | 90.0 | 81.5 | 84.3 | 53.0 | 74.4 | 79.8 | 84.0 | 46.5 | 50.2 | 78.6 | 72.2 | 88.4 | 86.9 | 88.1 | 53.7 | 75.8 | 82.1 | 90.0 | 44.7 | 48.7 | 73.9 | 73.2 |
| *KMC* | 93.7 | 87.8 | 91.7 | 57.2 | 76.6 | 87.8 | 92.9 | 59.6 | 56.0 | 84.4 | 78.8 | 87.7 | 82.1 | 89.7 | 54.5 | 68.4 | 81.0 | 48.8 | 56.4 | 52.8 | 79.9 | 70.1 | 88.0 | 84.3 | 93.2 | 58.1 | 76.6 | 85.8 | 57.2 | 49.0 | 48.1 | 81.8 | 72.2 |
| *MTLD* | 67.1 | 86.2 | 50.0 | 56.7 | 70.5 | 89.1 | 86.6 | 51.4 | 57.8 | 82.8 | 69.8 | 91.7 | 73.0 | 91.3 | 54.1 | 62.5 | 87.5 | 75.6 | 47.3 | 39.6 | 82.3 | 70.5 | 88.4 | 53.5 | 92.3 | 55.1 | 64.5 | 85.2 | 67.8 | 48.1 | 48.8 | 79.8 | 68.4 |
| *DSIR* | 90.0 | 6.9 | 90.8 | 47.4 | 70.3 | 90.5 | 54.6 | 49.5 | 37.7 | 83.1 | 62.1 | 85.6 | 77.2 | 84.8 | 47.3 | 73.1 | 81.5 | 59.4 | 37.1 | 55.0 | 70.7 | 67.2 | 84.6 | 80.7 | 84.1 | 52.0 | 72.2 | 81.4 | 74.7 | 48.2 | 50.2 | 70.1 | 69.8 |

Table 6: The detailed performance results of different language subset on *XStoryCloze* dataset.

Table 7 — Model: LLaMA-3.1-8B, Dataset: MMMLU

| Method | en | ar | bn | de | es | fr | hi | sw | yo | zh | Avg. | en | ar | bn | de | es | fr | hi | sw | yo | zh | Avg. | en | ar | bn | de | es | fr | hi | sw | yo | zh | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | 1% | | | | | | | | | | | 3% | | | | | | | | | | | | 5% | | | | | | | |
| *Rand* | 52.5 | 31.9 | 29.6 | 42.4 | 45.2 | 40.8 | 26.5 | 25.5 | 13.5 | 31.3 | 34.2 | 54.7 | 39.2 | 29.0 | 44.8 | 46.6 | 46.8 | 33.2 | 26.8 | 17.2 | 43.9 | 37.4 | 56.6 | 38.3 | 28.5 | 46.4 | 50.0 | 45.1 | 35.0 | 28.6 | 21.1 | 43.9 | 38.1 |
| *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC** | 58.4 | 40.8 | 22.1 | 45.9 | 48.9 | 41.3 | 35.3 | 31.9 | 13.4 | 43.0 | 38.1 | 56.5 | 38.9 | 29.2 | 47.3 | 52.4 | 46.2 | 33.4 | 26.9 | 22.7 | 42.3 | 39.6 | 57.7 | 40.1 | 27.2 | 44.5 | 47.7 | 39.2 | 34.6 | 21.8 | 13.7 | 45.3 | 37.2 |
| *MTLD* | 43.6 | 7.8 | 2.2 | 39.1 | 38.5 | 25.1 | 28.2 | 16.6 | 13.1 | 36.6 | 25.1 | 53.4 | 33.8 | 6.7 | 43.5 | 42.6 | 34.0 | 31.1 | 24.9 | 15.9 | 39.2 | 32.5 | 50.0 | 34.6 | 1.9 | 42.9 | 46.2 | 32.8 | 24.2 | 20.3 | 6.8 | 38.8 | 29.9 |
| *Nat* | 52.0 | 36.9 | 30.7 | 38.7 | 0.4 | 1.7 | 36.9 | 21.2 | 3.0 | 23.0 | 24.5 | 54.1 | 40.1 | 27.7 | 45.5 | 45.8 | 34.4 | 36.8 | 21.2 | 7.3 | 39.8 | 35.3 | 57.1 | 40.7 | 27.7 | 46.4 | 44.4 | 42.8 | 37.9 | 25.0 | 9.8 | 44.4 | 37.6 |
| *Coh** | 56.5 | 23.5 | 28.1 | 46.6 | 46.6 | 42.1 | 37.9 | 20.8 | 1.9 | 42.3 | 34.6 | 56.2 | 36.2 | 22.5 | 45.1 | 49.2 | 47.9 | 33.1 | 19.7 | 4.9 | 43.7 | 35.9 | 56.7 | 37.4 | 31.0 | 46.0 | 51.4 | 28.7 | 35.2 | 27.3 | 7.7 | 43.0 | 36.4 |
| *Und* | 54.7 | 37.2 | 28.5 | 36.1 | 4.8 | 20.8 | 35.9 | 25.9 | 5.3 | 39.0 | 28.8 | 51.4 | 41.3 | 24.8 | 46.2 | 45.2 | 37.3 | 26.8 | 27.8 | 6.1 | 43.2 | 35.0 | 56.8 | 42.0 | 30.6 | 47.0 | 47.9 | 45.5 | 34.4 | 25.3 | 9.2 | 45.7 | 38.4 |
| *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR* | 46.6 | 21.1 | 7.8 | 37.8 | 40.1 | 33.6 | 27.1 | 16.9 | 8.8 | 33.5 | 27.3 | 53.3 | 22.7 | 22.5 | 41.4 | 49.3 | 35.0 | 34.0 | 13.2 | 1.7 | 27.4 | 30.1 | 52.8 | 33.2 | 32.1 | 44.8 | 49.6 | 38.8 | 36.5 | 19.7 | 4.0 | 39.1 | 35.1 |
| *LESS** | 56.3 | 36.9 | 32.2 | 46.0 | 50.4 | 47.1 | 36.7 | 28.9 | 26.2 | 43.5 | 40.4 | 55.8 | 38.2 | 29.6 | 47.4 | 49.5 | 49.0 | 35.4 | 31.4 | 27.3 | 42.8 | 40.6 | 58.8 | 35.9 | 29.0 | 47.5 | 49.5 | 48.8 | 34.8 | 30.6 | 28.7 | 44.3 | 40.8 |
| *Applying LangGPS (to random selection and three best well-performing baselines marked with *)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 54.7 | 31.4 | 19.0 | 44.0 | 46.3 | 38.3 | 31.5 | 24.0 | 9.2 | 38.4 | 33.7 | 55.6 | 38.0 | 27.4 | 42.0 | 47.3 | 44.2 | 35.8 | 26.9 | 19.8 | 44.8 | 38.2 | 57.4 | 37.5 | 27.4 | 45.1 | 50.7 | 45.6 | 36.3 | 30.2 | 22.5 | 44.3 | 39.7 |
| *KMC* | 53.4 | 36.7 | 30.8 | 43.1 | 48.1 | 38.3 | 34.4 | 30.0 | 26.7 | 41.1 | 38.3 | 57.5 | 38.1 | 31.7 | 48.5 | 47.8 | 44.8 | 34.0 | 30.0 | 24.7 | 44.4 | 40.2 | 56.5 | 40.5 | 18.9 | 42.3 | 48.5 | 44.5 | 35.1 | 31.6 | 26.8 | 43.7 | 38.8 |
| *Coh* | 57.9 | 34.8 | 33.7 | 46.2 | 51.8 | 46.3 | 37.1 | 26.0 | 0.7 | 41.8 | 37.6 | 50.9 | 35.9 | 26.1 | 44.5 | 49.9 | 40.7 | 36.7 | 24.2 | 1.7 | 43.6 | 35.4 | 52.5 | 36.5 | 29.3 | 42.4 | 44.7 | 42.4 | 32.0 | 29.3 | 24.9 | 41.0 | 37.5 |
| *LESS* | 56.0 | 37.6 | 28.0 | 45.8 | 48.0 | 46.3 | 34.0 | 31.0 | 28.5 | 43.6 | 39.9 | 57.9 | 37.8 | 29.0 | 45.0 | 50.4 | 45.8 | 35.5 | 27.8 | 26.4 | 44.5 | 40.0 | 59.0 | 38.1 | 31.6 | 48.6 | 51.3 | 49.8 | 37.5 | 29.7 | 29.5 | 46.5 | 42.2 |

Model: Qwen2.5-7B, Dataset: MMMLU

| Method | en | ar | el | es | fr | hi | sw | tr | vi | zh | Avg. | en | ar | el | es | fr | hi | sw | tr | vi | zh | Avg. | en | ar | el | es | fr | hi | sw | tr | vi | zh | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | 1% | | | | | | | | | | | 3% | | | | | | | | | | | | 5% | | | | | | | |
| *Rand* | 68.1 | 52.1 | 37.4 | 60.3 | 64.2 | 61.5 | 39.0 | 26.8 | 21.6 | 62.2 | 49.3 | 55.9 | 53.1 | 32.7 | 60.2 | 62.9 | 62.6 | 36.3 | 26.6 | 19.5 | 46.2 | 45.6 | 53.6 | 52.4 | 34.3 | 59.2 | 59.7 | 62.7 | 38.4 | 25.6 | 16.2 | 39.4 | 44.1 |
| *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC** | 58.3 | 52.5 | 23.9 | 59.0 | 64.5 | 61.5 | 41.1 | 26.5 | 15.3 | 35.9 | 43.9 | 48.3 | 52.3 | 37.0 | 59.0 | 62.0 | 62.6 | 31.6 | 23.4 | 22.9 | 37.1 | 43.6 | 63.7 | 51.2 | 40.8 | 60.6 | 58.8 | 61.7 | 35.6 | 26.7 | 19.4 | 49.6 | 46.8 |
| *MTLD** | 64.7 | 52.4 | 35.2 | 61.0 | 62.5 | 62.6 | 43.2 | 18.3 | 7.2 | 57.7 | 46.5 | 67.3 | 48.6 | 29.1 | 58.7 | 64.6 | 62.7 | 31.9 | 23.1 | 13.6 | 61.0 | 46.1 | 65.2 | 51.1 | 23.8 | 58.1 | 61.4 | 60.5 | 34.2 | 20.5 | 3.6 | 58.1 | 43.7 |
| *Nat* | 66.9 | 50.7 | 38.5 | 58.7 | 64.3 | 60.8 | 33.6 | 28.8 | 1.2 | 61.3 | 46.5 | 68.1 | 53.3 | 25.2 | 58.2 | 61.2 | 63.2 | 34.0 | 18.6 | 8.0 | 51.4 | 44.1 | 58.4 | 53.4 | 23.4 | 60.0 | 62.4 | 61.5 | 26.4 | 20.6 | 16.6 | 42.2 | 42.5 |
| *Coh* | 64.9 | 46.7 | 27.6 | 58.6 | 62.2 | 59.1 | 14.6 | 10.7 | 0.4 | 62.0 | 40.7 | 63.0 | 52.0 | 34.4 | 59.7 | 64.6 | 61.0 | 30.6 | 17.2 | 2.1 | 59.9 | 44.5 | 63.3 | 39.6 | 23.7 | 58.3 | 54.8 | 56.4 | 17.5 | 12.2 | 2.1 | 52.2 | 38.0 |
| *Und* | 68.1 | 49.7 | 30.6 | 59.7 | 64.4 | 61.6 | 29.3 | 10.0 | 1.6 | 63.0 | 43.8 | 68.1 | 54.1 | 25.3 | 59.5 | 64.2 | 62.5 | 30.3 | 26.3 | 5.9 | 54.7 | 45.1 | 56.5 | 54.0 | 24.9 | 58.9 | 56.1 | 59.9 | 27.9 | 27.6 | 14.4 | 41.3 | 42.2 |
| *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR** | 66.3 | 49.8 | 34.9 | 59.7 | 62.6 | 61.1 | 28.8 | 23.7 | 11.2 | 60.9 | 45.9 | 66.4 | 51.0 | 35.4 | 60.0 | 63.8 | 60.3 | 33.0 | 19.6 | 6.5 | 53.9 | 45.0 | 63.5 | 51.0 | 36.1 | 60.3 | 65.5 | 59.4 | 33.1 | 19.9 | 7.0 | 57.9 | 45.4 |
| *LESS* | 60.8 | 49.3 | 21.7 | 52.2 | 50.7 | 57.4 | 28.0 | 28.5 | 26.0 | 54.1 | 42.9 | 54.6 | 49.8 | 34.5 | 51.1 | 52.9 | 52.3 | 34.7 | 20.7 | 12.3 | 47.9 | 41.1 | 50.7 | 50.4 | 37.2 | 51.6 | 58.7 | 59.0 | 30.2 | 32.8 | 29.8 | 51.5 | 45.2 |
| *Applying LangGPS (to random selection and three best well-performing baselines marked with *)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 67.9 | 52.1 | 33.8 | 60.7 | 63.9 | 62.7 | 36.1 | 25.5 | 21.3 | 63.8 | 48.8 | 59.8 | 53.1 | 33.8 | 60.7 | 64.7 | 61.8 | 36.2 | 26.5 | 17.2 | 47.6 | 46.2 | 52.5 | 52.9 | 38.0 | 56.7 | 61.9 | 58.7 | 36.8 | 24.9 | 18.1 | 41.6 | 44.2 |
| *KMC* | 50.2 | 52.4 | 31.4 | 61.4 | 62.5 | 62.5 | 37.9 | 25.0 | 19.1 | 36.6 | 43.9 | 65.7 | 53.3 | 41.3 | 58.1 | 50.2 | 62.4 | 43.7 | 27.7 | 21.7 | 41.2 | 46.5 | 46.8 | 52.8 | 42.3 | 60.3 | 56.5 | 61.9 | 43.0 | 26.0 | 20.5 | 41.8 | 45.2 |
| *MTLD* | 66.8 | 50.5 | 36.6 | 60.0 | 57.9 | 60.6 | 39.7 | 25.6 | 14.8 | 59.9 | 47.2 | 66.9 | 53.4 | 41.0 | 62.0 | 63.8 | 64.4 | 41.3 | 23.2 | 3.9 | 56.2 | 47.6 | 56.0 | 53.7 | 39.9 | 60.2 | 59.6 | 64.0 | 37.9 | 21.4 | 0.9 | 39.4 | 43.3 |
| *DSIR* | 66.9 | 53.8 | 38.7 | 61.7 | 64.7 | 63.9 | 43.6 | 26.4 | 7.0 | 62.4 | 48.9 | 68.9 | 54.6 | 41.0 | 61.1 | 64.7 | 63.9 | 38.4 | 28.0 | 16.4 | 64.8 | 50.2 | 59.40 | 54.50 | 38.30 | 60.80 | 64.90 | 62.10 | 42.50 | 26.40 | 8.60 | 51.50 | 46.9 |

Table 7: The detailed performance results of different language subset on *MMMLU* dataset.

| | Model: LLaMA-3.1-8B | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Method** | Dataset: MKQA | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1% | | | | | | | | | | | 3% | | | | | | | | | | | 5% | | | | | | | | | | |
| Langs | en | ar | de | ja | ko | pt | ru | tr | vi | zh | **Avg.** | en | ar | de | ja | ko | pt | ru | tr | vi | zh | **Avg.** | en | ar | de | ja | ko | pt | ru | tr | vi | zh | **Avg.** |
| *Rand* | 31.2 | 7.0 | 21.4 | 10.0 | 6.1 | 18.5 | 10.8 | 16.2 | 16.6 | 8.6 | 14.6 | 31.9 | 5.7 | 22.7 | 8.8 | 5.4 | 19.3 | 8.9 | 15.9 | 16.4 | 9.1 | 14.4 | 30.9 | 6.0 | 22.1 | 8.3 | 5.0 | 20.2 | 9.9 | 15.7 | 17.2 | 9.1 | 14.4 |
| | *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC** | 31.4 | 7.4 | 21.2 | 11.0 | 6.0 | 18.9 | 10.6 | 16.7 | 19.0 | 8.6 | 15.1 | 33.5 | 5.9 | 22.8 | 8.0 | 5.4 | 20.6 | 9.3 | 15.5 | 17.2 | 9.7 | 14.8 | 32.0 | 6.6 | 22.2 | 8.2 | 6.0 | 20.8 | 9.9 | 16.7 | 17.2 | 10.4 | 15.0 |
| *MTLD* | 30.2 | 8.5 | 21.6 | 11.4 | 6.5 | 20.6 | 11.0 | 18.7 | 18.9 | 9.6 | 15.7 | 33.6 | 8.1 | 22.2 | 11.7 | 7.0 | 21.4 | 9.1 | 17.7 | 17.8 | 10.2 | 15.9 | 29.1 | 7.8 | 19.2 | 10.4 | 7.7 | 18.2 | 9.6 | 18.6 | 16.1 | 8.8 | 14.5 |
| *Nat* | 37.2 | 7.7 | 25.8 | 13.2 | 8.0 | 26.6 | 12.1 | 20.4 | 22.6 | 11.1 | 18.5 | 35.8 | 7.1 | 24.3 | 13.8 | 7.2 | 23.2 | 11.2 | 19.2 | 19.3 | 12.0 | 17.3 | 36.1 | 7.3 | 24.1 | 12.1 | 7.9 | 25.1 | 11.9 | 17.1 | 17.7 | 11.2 | 17.1 |
| *Coh** | 37.0 | 7.2 | 25.2 | 14.4 | 8.5 | 25.5 | 12.9 | 21.1 | 19.2 | 11.0 | 18.2 | 36.1 | 8.1 | 26.6 | 13.8 | 8.9 | 25.2 | 11.2 | 20.1 | 19.2 | 11.6 | 18.1 | 37.1 | 7.6 | 26.1 | 14.4 | 8.8 | 25.7 | 11.7 | 20.9 | 20.6 | 10.8 | 18.4 |
| *Und* | 36.9 | 8.2 | 25.0 | 13.9 | 8.9 | 25.6 | 11.8 | 19.0 | 21.3 | 11.6 | 18.2 | 35.3 | 5.1 | 24.8 | 12.7 | 8.6 | 23.8 | 11.0 | 18.1 | 18.2 | 11.5 | 16.9 | 35.3 | 7.5 | 22.7 | 12.0 | 8.7 | 25.9 | 12.5 | 20.1 | 19.2 | 10.0 | 17.4 |
| | *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR* | 35.4 | 8.9 | 25.4 | 12.9 | 8.8 | 25.1 | 12.2 | 19.3 | 18.8 | 9.9 | 17.7 | 35.0 | 7.8 | 22.5 | 12.0 | 9.2 | 22.4 | 11.1 | 19.5 | 19.1 | 10.2 | 16.9 | 32.0 | 7.2 | 23.2 | 12.1 | 8.7 | 22.4 | 10.3 | 19.3 | 18.4 | 10.0 | 16.4 |
| *LESS** | 36.0 | 8.9 | 25.1 | 12.5 | 8.5 | 25.4 | 11.6 | 21.9 | 20.1 | 11.5 | 18.2 | 35.3 | 8.0 | 25.4 | 12.7 | 7.0 | 26.8 | 13.2 | 20.4 | 20.5 | 10.9 | 18.0 | 34.7 | 7.5 | 24.9 | 11.9 | 6.9 | 26.0 | 11.1 | 19.5 | 19.0 | 10.6 | 17.2 |
| | *Applying LangGPS (to random selection and three best well-performing baselines marked with *)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 33.0 | 7.7 | 21.4 | 10.0 | 6.4 | 21.1 | 11.0 | 16.9 | 17.4 | 8.8 | 15.4 | 32.5 | 6.4 | 22.9 | 9.4 | 6.1 | 21.3 | 9.7 | 16.9 | 16.1 | 8.3 | 15.0 | 33.6 | 6.4 | 22.2 | 9.1 | 5.5 | 19.6 | 9.7 | 16.2 | 16.1 | 8.7 | 14.7 |
| *KMC* | 32.9 | 8.0 | 22.3 | 9.4 | 5.8 | 21.9 | 9.6 | 17.3 | 17.3 | 8.4 | 15.3 | 31.9 | 6.9 | 21.3 | 10.6 | 6.4 | 20.1 | 10.7 | 15.7 | 17.2 | 9.8 | 15.1 | 29.6 | 6.3 | 19.7 | 10.9 | 4.9 | 17.2 | 9.5 | 16.1 | 14.2 | 9.3 | 13.8 |
| *Coh* | 37.6 | 8.3 | 26.0 | 14.4 | 8.8 | 26.3 | 14.0 | 21.6 | 22.0 | 11.7 | 19.1 | 36.3 | 7.8 | 24.9 | 14.0 | 9.3 | 26.8 | 10.7 | 18.0 | 22.1 | 11.6 | 18.2 | 35.2 | 7.8 | 24.6 | 12.8 | 9.1 | 26.2 | 11.0 | 19.6 | 18.7 | 11.0 | 17.6 |
| *LESS* | 34.5 | 8.1 | 24.0 | 11.4 | 7.7 | 23.6 | 11.6 | 18.8 | 20.2 | 10.8 | 17.1 | 33.9 | 7.1 | 24.1 | 11.0 | 7.9 | 23.1 | 10.1 | 18.5 | 18.7 | 10.5 | 16.5 | 35.2 | 7.9 | 24.9 | 11.5 | 7.5 | 24.6 | 11.7 | 18.4 | 18.6 | 10.1 | 17.0 |

| | Model: Qwen2.5-7B | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Method** | Dataset: MKQA | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1% | | | | | | | | | | | 3% | | | | | | | | | | | 5% | | | | | | | | | | |
| Langs | en | ar | de | ja | ko | pt | ru | tr | vi | zh | **Avg.** | en | ar | de | ja | ko | pt | ru | tr | vi | zh | **Avg.** | en | ar | de | ja | ko | pt | ru | tr | vi | zh | **Avg.** |
| *Rand* | 31.9 | 6.9 | 19.4 | 11.9 | 8.3 | 21.6 | 10.2 | 12.8 | 16.4 | 15.7 | 15.5 | 30.5 | 6.7 | 19.3 | 11.9 | 7.9 | 20.0 | 9.5 | 13.1 | 16.6 | 15.2 | 15.1 | 28.6 | 6.5 | 19.0 | 10.8 | 7.7 | 20.4 | 9.0 | 13.5 | 15.7 | 14.3 | 14.5 |
| | *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC** | 28.4 | 6.1 | 17.0 | 9.7 | 7.2 | 17.6 | 6.6 | 9.8 | 14.6 | 14.2 | 13.1 | 29.4 | 6.0 | 17.1 | 9.1 | 7.4 | 20.2 | 10.4 | 12.4 | 15.3 | 14.7 | 14.2 | 28.2 | 5.5 | 16.3 | 10.1 | 6.7 | 19.1 | 7.9 | 11.8 | 14.5 | 13.8 | 13.4 |
| *MTLD** | 32.3 | 7.4 | 19.0 | 12.4 | 8.1 | 20.9 | 9.7 | 12.0 | 15.2 | 15.2 | 15.2 | 29.7 | 6.7 | 17.6 | 12.6 | 7.7 | 18.7 | 8.5 | 12.2 | 13.6 | 15.2 | 14.2 | 29.3 | 7.3 | 18.5 | 11.8 | 7.7 | 20.1 | 9.0 | 12.7 | 14.8 | 14.6 | 14.6 |
| *Nat* | 31.2 | 6.3 | 19.7 | 11.8 | 9.1 | 20.3 | 8.9 | 13.4 | 16.6 | 15.0 | 15.2 | 30.4 | 6.5 | 17.7 | 12.7 | 7.7 | 18.9 | 8.5 | 11.8 | 15.6 | 15.3 | 14.5 | 29.2 | 6.6 | 19.0 | 11.4 | 7.5 | 17.7 | 8.6 | 12.1 | 16.2 | 15.5 | 14.4 |
| *Coh* | 30.1 | 6.5 | 18.4 | 12.3 | 8.9 | 20.4 | 9.4 | 12.8 | 18.2 | 15.6 | 15.3 | 28.1 | 6.2 | 17.6 | 12.4 | 8.3 | 20.1 | 9.5 | 13.1 | 16.7 | 14.6 | 14.7 | 28.8 | 5.8 | 15.2 | 12.9 | 9.4 | 19.0 | 10.1 | 10.9 | 17.0 | 15.9 | 14.5 |
| *Und* | 30.8 | 6.5 | 18.7 | 12.7 | 9.1 | 19.2 | 8.9 | 12.1 | 16.7 | 14.7 | 14.9 | 29.2 | 6.8 | 18.0 | 12.4 | 8.7 | 18.5 | 8.6 | 11.9 | 15.4 | 14.5 | 14.4 | 29.1 | 7.5 | 18.3 | 11.1 | 8.0 | 19.7 | 9.2 | 13.5 | 16.0 | 14.4 | 14.7 |
| | *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR** | 31.1 | 7.3 | 21.2 | 13.1 | 8.3 | 23.6 | 10.2 | 14.7 | 18.8 | 16.0 | 16.4 | 27.6 | 5.5 | 17.2 | 11.2 | 7.9 | 17.7 | 8.7 | 11.7 | 13.6 | 14.5 | 13.6 | 27.7 | 5.7 | 17.7 | 11.0 | 8.5 | 20.1 | 8.4 | 12.3 | 15.7 | 14.8 | 14.2 |
| *LESS* | 23.7 | 5.8 | 18.5 | 10.1 | 7.1 | 18.2 | 9.1 | 14.3 | 16.3 | 15.8 | 13.9 | 25.1 | 6.3 | 19.0 | 9.5 | 7.0 | 20.6 | 9.5 | 14.9 | 16.8 | 15.6 | 14.4 | 21.4 | 5.8 | 17.0 | 9.3 | 7.3 | 16.7 | 8.6 | 13.2 | 14.4 | 14.7 | 12.8 |
| | *Applying LangGPS (to random selection and three best well-performing baselines marked with *)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 33.1 | 6.6 | 19.1 | 11.9 | 8.1 | 21.8 | 9.6 | 13.0 | 15.8 | 15.2 | 15.4 | 30.7 | 6.6 | 18.8 | 11.4 | 8.0 | 20.2 | 9.2 | 13.2 | 14.8 | 15.1 | 14.8 | 29.5 | 6.5 | 19.6 | 10.7 | 7.2 | 20.3 | 8.8 | 12.8 | 15.4 | 14.5 | 14.5 |
| *KMC* | 28.1 | 6.1 | 20.4 | 9.9 | 6.9 | 17.8 | 9.4 | 11.1 | 13.3 | 14.6 | 13.8 | 29.6 | 6.9 | 18.3 | 10.1 | 7.7 | 20.1 | 9.4 | 12.3 | 15.6 | 13.6 | 14.4 | 28.9 | 5.7 | 19.7 | 10.4 | 6.6 | 18.8 | 9.1 | 11.8 | 15.3 | 13.1 | 13.9 |
| *MTLD* | 30.9 | 7.3 | 19.7 | 12.4 | 8.6 | 21.9 | 9.3 | 12.3 | 14.6 | 14.5 | 15.2 | 29.6 | 7.0 | 18.6 | 11.5 | 8.9 | 21.3 | 8.4 | 13.1 | 16.5 | 14.6 | 15.0 | 29.1 | 5.9 | 19.2 | 12.3 | 7.9 | 20.9 | 9.2 | 12.9 | 16.3 | 14.9 | 14.9 |
| *DSIR* | 31.0 | 7.3 | 20.2 | 12.8 | 9.5 | 21.1 | 9.4 | 13.6 | 18.0 | 15.9 | 15.9 | 29.2 | 6.3 | 18.3 | 11.9 | 8.3 | 19.8 | 8.5 | 13.6 | 17.6 | 14.8 | 14.8 | 8.8 | 6.9 | 18.7 | 11.8 | 9.2 | 19.9 | 8.6 | 13.2 | 16.9 | 15.8 | 15.0 |

Table 8: The detailed performance results of different language subset on *MKQA* dataset.

| | Model: LLaMA-3.1-8B | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Method** | Dataset: XQuAD | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1% | | | | | | | | | | | 3% | | | | | | | | | | | 5% | | | | | | | | | | | |
| Langs | en | ar | de | el | es | hi | ru | th | tr | zh | **Avg.** | en | ar | de | el | es | hi | ru | th | tr | zh | **Avg.** | en | ar | de | el | es | hi | ru | th | tr | zh | **Avg.** |
| *Rand* | 75.1 | 58.4 | 67.6 | 53.9 | 68.3 | 62.5 | 53.8 | 63.7 | 58.0 | 74.9 | 63.6 | 73.0 | 52.9 | 64.5 | 53.2 | 66.6 | 59.2 | 50.1 | 55.5 | 52.4 | 72.2 | 59.9 | 69.7 | 53.9 | 63.8 | 50.9 | 64.5 | 60.9 | 48.6 | 60.4 | 51.8 | 71.9 | 59.6 |
| | *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC\** | 71.6 | 54.2 | 63.1 | 51.8 | 60.9 | 62.9 | 48.7 | 61.7 | 48.5 | 73.3 | 59.7 | 72.7 | 52.9 | 61.9 | 50.0 | 65.1 | 61.8 | 47.8 | 56.6 | 49.6 | 70.3 | 58.9 | 72.7 | 53.6 | 61.9 | 50.8 | 66.5 | 60.7 | 51.0 | 57.7 | 47.4 | 71.0 | 59.3 |
| *MTLD* | 74.3 | 55.5 | 64.8 | 52.5 | 62.7 | 60.4 | 50.5 | 59.2 | 51.8 | 71.4 | 60.3 | 78.7 | 55.5 | 65.7 | 54.9 | 66.1 | 63.5 | 53.4 | 59.6 | 51.9 | 75.5 | 62.5 | 65.3 | 52.9 | 61.8 | 50.8 | 63.8 | 59.7 | 47.7 | 56.0 | 46.4 | 67.6 | 57.2 |
| *Nat* | 79.0 | 51.8 | 63.1 | 52.4 | 65.9 | 54.1 | 48.5 | 50.2 | 56.1 | 71.7 | 59.3 | 75.2 | 57.2 | 63.6 | 52.9 | 65.2 | 60.2 | 53.9 | 54.9 | 56.1 | 74.3 | 61.3 | 74.3 | 54.5 | 59.7 | 51.4 | 64.9 | 58.0 | 51.7 | 52.9 | 55.0 | 72.6 | 59.5 |
| *Coh\** | 79.2 | 58.3 | 67.4 | 57.1 | 71.2 | 58.2 | 56.0 | 54.5 | 62.8 | 76.6 | 64.1 | 79.2 | 59.2 | 66.3 | 56.6 | 73.1 | 59.3 | 55.4 | 57.8 | 60.4 | 78.2 | 64.6 | 82.0 | 58.3 | 65.8 | 56.9 | 73.2 | 58.0 | 56.6 | 56.7 | 60.8 | 79.3 | 64.8 |
| *Und* | 78.2 | 55.0 | 61.8 | 52.6 | 64.6 | 58.2 | 52.2 | 52.1 | 59.3 | 71.8 | 60.6 | 75.5 | 55.8 | 60.8 | 52.9 | 66.6 | 57.2 | 52.1 | 54.1 | 56.6 | 74.6 | 60.6 | 74.3 | 55.9 | 60.8 | 52.1 | 65.2 | 55.7 | 51.9 | 50.8 | 55.2 | 72.9 | 59.5 |
| | *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR* | 80.9 | 59.8 | 68.2 | 55.2 | 69.1 | 56.6 | 55.0 | 61.7 | 64.5 | 76.1 | 64.7 | 74.2 | 57.6 | 65.2 | 56.1 | 69.3 | 60.8 | 57.3 | 60.6 | 56.7 | 76.1 | 63.4 | 72.0 | 56.6 | 65.6 | 57.0 | 70.3 | 59.7 | 54.6 | 57.8 | 54.9 | 75.6 | 62.4 |
| *LESS\** | 80.7 | 58.4 | 67.5 | 57.2 | 67.7 | 63.1 | 56.5 | 62.8 | 60.8 | 77.5 | 65.2 | 79.6 | 59.3 | 69.3 | 56.8 | 69.0 | 61.8 | 58.5 | 63.3 | 61.4 | 77.0 | 65.6 | 78.7 | 59.7 | 65.7 | 58.4 | 69.4 | 61.3 | 59.3 | 62.4 | 56.9 | 76.5 | 64.8 |
| | *Applying LangGPS (to random selection and three best well-performing baselines marked with \*)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 75.1 | 58.4 | 64.6 | 55.6 | 67.3 | 60.5 | 52.4 | 62.1 | 55.6 | 74.5 | 62.6 | 74.3 | 54.0 | 63.3 | 51.8 | 63.0 | 59.9 | 51.0 | 58.8 | 50.2 | 72.1 | 60.2 | 73.9 | 56.0 | 65.5 | 53.3 | 68.3 | 62.1 | 52.0 | 60.5 | 51.1 | 74.1 | 62.0 |
| *KMC* | 74.3 | 56.2 | 66.1 | 51.8 | 68.7 | 62.9 | 50.9 | 59.6 | 51.7 | 75.5 | 61.8 | 73.9 | 54.6 | 63.4 | 53.3 | 65.9 | 61.0 | 51.6 | 58.7 | 50.0 | 70.3 | 60.3 | 70.8 | 55.8 | 64.6 | 53.4 | 66.3 | 63.0 | 52.6 | 57.1 | 50.3 | 71.6 | 60.6 |
| *Coh* | 80.4 | 58.2 | 62.8 | 56.1 | 69.3 | 58.7 | 55.9 | 53.8 | 62.8 | 75.8 | 63.4 | 78.5 | 54.9 | 63.8 | 56.0 | 69.4 | 56.7 | 55.1 | 53.9 | 62.3 | 75.1 | 62.6 | 77.2 | 59.3 | 65.3 | 56.7 | 72.7 | 59.1 | 57.5 | 56.5 | 59.1 | 73.9 | 63.7 |
| *LESS* | 79.4 | 60.1 | 67.2 | 58.8 | 69.1 | 63.6 | 56.6 | 63.2 | 60.8 | 76.6 | 65.6 | 78.9 | 59.7 | 67.0 | 56.2 | 69.8 | 61.5 | 56.6 | 63.0 | 56.1 | 74.4 | 64.3 | 78.8 | 60.3 | 64.6 | 56.6 | 69.2 | 60.8 | 58.2 | 60.4 | 58.7 | 75.1 | 64.3 |
| | Model: Qwen2.5-7B | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| **Method** | Dataset: XQuAD | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1% | | | | | | | | | | | 3% | | | | | | | | | | | 5% | | | | | | | | | | | |
| Langs | en | ar | de | el | es | hi | ru | th | tr | zh | **Avg.** | en | ar | de | el | es | hi | ru | th | tr | zh | **Avg.** | en | ar | de | el | es | hi | ru | th | tr | zh | **Avg.** |
| *Rand* | 77.6 | 65.4 | 69.9 | 46.0 | 71.2 | 54.4 | 52.7 | 74.4 | 58.9 | 83.7 | 65.4 | 77.5 | 64.9 | 70.3 | 42.1 | 72.2 | 47.7 | 53.8 | 69.9 | 57.7 | 83.5 | 64.0 | 78.3 | 64.8 | 70.4 | 40.6 | 73.9 | 47.5 | 54.4 | 67.1 | 58.4 | 83.7 | 63.9 |
| | *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC\** | 73.7 | 59.8 | 68.0 | 38.9 | 66.5 | 52.3 | 47.1 | 70.8 | 54.0 | 81.3 | 61.2 | 78.1 | 66.1 | 68.6 | 36.1 | 72.0 | 44.2 | 52.3 | 65.9 | 53.3 | 83.9 | 62.0 | 78.1 | 63.4 | 71.1 | 39.4 | 71.3 | 44.4 | 52.1 | 65.3 | 56.1 | 82.9 | 62.4 |
| *MTLD\** | 78.6 | 65.0 | 71.5 | 44.5 | 72.1 | 54.1 | 53.7 | 73.2 | 56.0 | 86.8 | 65.5 | 78.6 | 61.3 | 70.0 | 40.8 | 71.4 | 51.6 | 55.9 | 67.2 | 53.2 | 85.2 | 63.5 | 79.5 | 62.1 | 70.4 | 38.0 | 73.3 | 46.8 | 53.5 | 65.3 | 53.9 | 84.8 | 62.8 |
| *Nat* | 79.2 | 64.1 | 70.3 | 21.5 | 71.2 | 35.7 | 53.2 | 52.9 | 61.4 | 81.9 | 59.2 | 80.9 | 66.7 | 72.0 | 27.2 | 72.4 | 40.3 | 57.7 | 59.6 | 59.9 | 84.8 | 62.2 | 79.3 | 63.4 | 69.5 | 32.9 | 70.7 | 40.0 | 55.9 | 60.5 | 56.6 | 86.3 | 61.5 |
| *Coh* | 82.2 | 67.7 | 73.1 | 17.7 | 76.1 | 28.0 | 56.9 | 50.8 | 63.4 | 86.5 | 60.2 | 82.8 | 67.0 | 73.2 | 21.3 | 77.4 | 33.8 | 59.2 | 54.5 | 61.3 | 87.4 | 61.8 | 81.3 | 65.0 | 72.3 | 21.1 | 76.1 | 32.8 | 58.9 | 54.0 | 62.4 | 87.2 | 61.1 |
| *Und* | 77.7 | 61.8 | 69.9 | 21.4 | 71.0 | 34.4 | 54.1 | 52.4 | 59.3 | 81.2 | 58.3 | 80.0 | 66.0 | 70.7 | 26.1 | 73.7 | 38.9 | 58.2 | 58.5 | 61.1 | 85.2 | 61.8 | 78.8 | 61.9 | 69.2 | 36.1 | 70.4 | 41.8 | 57.0 | 62.7 | 58.2 | 86.1 | 62.2 |
| | *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR\** | 82.9 | 63.9 | 72.0 | 27.5 | 76.6 | 30.1 | 59.1 | 59.9 | 63.9 | 85.5 | 62.1 | 79.8 | 62.8 | 71.7 | 30.7 | 72.9 | 35.6 | 55.5 | 64.2 | 58.4 | 84.5 | 61.6 | 80.3 | 62.9 | 70.8 | 34.6 | 74.9 | 35.8 | 55.1 | 61.7 | 57.4 | 83.9 | 61.7 |
| *LESS* | 73.7 | 62.1 | 68.1 | 51.1 | 61.3 | 63.9 | 51.3 | 77.2 | 55.4 | 85.5 | 65.0 | 69.7 | 63.3 | 68.2 | 46.5 | 63.9 | 65.1 | 51.8 | 76.0 | 51.9 | 85.5 | 64.2 | 77.7 | 66.0 | 71.9 | 50.5 | 68.7 | 65.2 | 57.9 | 77.1 | 57.6 | 86.8 | 68.0 |
| | *Applying LangGPS (to random selection and three best well-performing baselines marked with \*)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 77.9 | 64.8 | 69.2 | 45.4 | 70.5 | 50.6 | 52.0 | 73.7 | 59.0 | 84.2 | 64.7 | 79.0 | 64.7 | 71.4 | 43.2 | 73.1 | 45.5 | 55.8 | 71.1 | 61.5 | 85.1 | 65.0 | 79.1 | 65.7 | 72.2 | 39.8 | 74.5 | 47.6 | 55.4 | 68.2 | 62.3 | 84.6 | 64.9 |
| *KMC* | 77.0 | 61.1 | 69.3 | 41.8 | 71.3 | 49.6 | 50.5 | 70.3 | 56.3 | 79.6 | 62.7 | 78.2 | 62.6 | 69.2 | 35.0 | 71.0 | 46.1 | 51.3 | 65.8 | 59.3 | 84.9 | 62.3 | 79.7 | 62.6 | 70.8 | 40.1 | 73.1 | 47.8 | 55.0 | 64.4 | 57.3 | 83.0 | 63.4 |
| *MTLD* | 81.3 | 64.5 | 71.3 | 44.0 | 72.9 | 53.1 | 55.0 | 70.7 | 57.6 | 84.0 | 65.5 | 80.4 | 66.0 | 70.8 | 40.7 | 74.1 | 46.9 | 55.3 | 66.2 | 57.7 | 84.9 | 64.3 | 82.0 | 65.4 | 72.2 | 35.3 | 76.1 | 42.4 | 57.1 | 62.1 | 60.3 | 86.2 | 63.9 |
| *DSIR* | 81.6 | 64.7 | 72.7 | 33.7 | 76.1 | 42.8 | 58.5 | 63.9 | 63.1 | 84.5 | 64.2 | 81.0 | 65.5 | 73.2 | 32.0 | 75.5 | 43.1 | 56.8 | 61.9 | 63.6 | 84.5 | 63.7 | 80.8 | 64.0 | 73.4 | 33.9 | 75.5 | 46.3 | 55.5 | 61.8 | 63.8 | 85.7 | 64.1 |

Table 9: The detailed performance results of different language subset on *XQuAD* dataset.

| | Model: LLaMA-3.1-8B | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Method** | **Dataset: XLSum** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1% | | | | | | | | | | | 3% | | | | | | | | | | | 5% | | | | | | | | | | |
| Langs | en | es | fr | ko | pt | sw | tr | uk | vi | zh | **Avg.** | en | es | fr | ko | pt | sw | tr | uk | vi | zh | **Avg.** | en | es | fr | ko | pt | sw | tr | uk | vi | zh | **Avg.** |
| *Rand* | 26.1 | 21.8 | 27.3 | 14.5 | 24.2 | 22.6 | 20.2 | 17.5 | 18.9 | 25.9 | 21.9 | 29.3 | 21.6 | 28.9 | 14.9 | 25.2 | 20.8 | 21.8 | 17.4 | 20.6 | 27.9 | 22.8 | 28.3 | 21.3 | 27.3 | 14.4 | 24.8 | 20.8 | 22.8 | 17.3 | 22.0 | 28.9 | 22.8 |
| | *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC\** | 26.3 | 20.1 | 24.7 | 15.3 | 23.0 | 18.7 | 19.7 | 17.2 | 19.2 | 25.8 | 21.0 | 29.3 | 22.0 | 27.8 | 16.1 | 25.5 | 21.2 | 22.4 | 15.8 | 22.1 | 29.1 | 23.1 | 27.1 | 21.5 | 28.2 | 13.8 | 23.4 | 22.0 | 23.1 | 17.0 | 20.6 | 27.7 | 22.4 |
| *MTLD* | 24.2 | 19.5 | 24.2 | 13.4 | 21.6 | 18.3 | 18.5 | 13.3 | 16.4 | 25.2 | 19.5 | 26.5 | 20.3 | 27.6 | 15.2 | 23.9 | 21.3 | 21.6 | 16.3 | 20.2 | 26.3 | 21.9 | 33.0 | 20.8 | 29.5 | 15.2 | 26.4 | 22.7 | 23.4 | 19.2 | 23.9 | 27.6 | 24.2 |
| *Nat* | 24.8 | 20.9 | 25.1 | 12.2 | 22.8 | 14.7 | 19.7 | 13.4 | 17.9 | 24.2 | 19.6 | 28.1 | 22.2 | 26.8 | 13.2 | 23.6 | 17.8 | 18.9 | 13.7 | 17.8 | 26.5 | 20.9 | 28.6 | 21.3 | 26.7 | 14.6 | 22.7 | 21.6 | 21.4 | 18.3 | 19.8 | 26.2 | 22.1 |
| *Coh\** | 26.1 | 21.3 | 26.6 | 15.1 | 25.4 | 23.2 | 22.7 | 18.2 | 19.2 | 25.8 | 22.4 | 30.6 | 21.0 | 28.1 | 15.6 | 26.1 | 23.6 | 24.3 | 19.5 | 23.2 | 28.4 | 24.0 | 29.1 | 21.5 | 28.5 | 15.0 | 26.4 | 23.8 | 24.3 | 17.2 | 23.1 | 28.6 | 23.8 |
| *Und* | 24.8 | 20.8 | 25.5 | 13.4 | 24.0 | 18.7 | 21.8 | 13.5 | 17.4 | 26.1 | 20.6 | 25.7 | 21.0 | 27.3 | 14.6 | 23.2 | 19.0 | 20.0 | 16.0 | 18.5 | 26.6 | 21.2 | 28.5 | 21.3 | 26.8 | 13.2 | 23.6 | 21.8 | 21.9 | 16.7 | 21.0 | 27.4 | 22.2 |
| | *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR* | 24.0 | 21.6 | 26.5 | 15.0 | 23.4 | 24.5 | 21.5 | 16.9 | 21.0 | 26.8 | 22.1 | 27.8 | 22.2 | 28.4 | 14.8 | 27.3 | 20.4 | 22.9 | 18.6 | 21.6 | 27.1 | 23.1 | 30.1 | 21.9 | 29.2 | 14.6 | 25.5 | 23.6 | 23.2 | 18.3 | 20.7 | 27.8 | 23.5 |
| *LESS\** | 24.9 | 20.7 | 24.8 | 14.5 | 22.4 | 22.6 | 20.6 | 16.3 | 19.3 | 26.2 | 21.2 | 25.0 | 20.8 | 25.9 | 15.1 | 23.4 | 21.6 | 21.3 | 16.7 | 19.4 | 27.4 | 21.7 | 27.1 | 21.4 | 27.6 | 15.9 | 22.9 | 22.6 | 22.5 | 17.5 | 22.0 | 28.9 | 22.8 |
| | *Applying LangGPS (to random selection and three best well-performing baselines marked with \*)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 26.5 | 21.3 | 26.7 | 14.2 | 24.7 | 21.4 | 21.9 | 17.4 | 19.1 | 26.0 | 21.9 | 29.1 | 21.6 | 27.9 | 15.8 | 23.8 | 20.1 | 24.1 | 18.0 | 20.3 | 27.5 | 22.8 | 30.6 | 22.1 | 27.1 | 15.9 | 25.7 | 23.7 | 24.1 | 18.3 | 22.3 | 28.4 | 23.8 |
| *KMC* | 24.9 | 20.8 | 25.3 | 14.5 | 22.8 | 21.9 | 21.5 | 17.9 | 20.6 | 24.9 | 21.5 | 29.8 | 21.5 | 27.1 | 14.3 | 24.2 | 18.6 | 22.7 | 15.0 | 20.5 | 28.0 | 22.2 | 30.2 | 21.2 | 28.4 | 16.2 | 24.0 | 24.5 | 24.9 | 19.1 | 21.8 | 28.2 | 23.8 |
| *Coh* | 25.5 | 21.2 | 27.0 | 13.6 | 25.2 | 23.6 | 22.0 | 17.4 | 21.4 | 25.0 | 22.2 | 29.9 | 21.4 | 29.2 | 16.7 | 27.0 | 25.0 | 24.3 | 19.0 | 22.6 | 28.2 | 24.3 | 32.1 | 22.4 | 28.6 | 16.8 | 24.8 | 24.5 | 22.9 | 18.4 | 23.4 | 29.1 | 24.3 |
| *LESS* | 26.8 | 21.2 | 26.5 | 16.5 | 23.7 | 23.7 | 20.8 | 17.5 | 21.0 | 26.9 | 22.5 | 30.5 | 20.8 | 29.2 | 15.0 | 20.7 | 21.1 | 23.5 | 17.2 | 22.5 | 27.6 | 22.8 | 30.6 | 20.7 | 29.0 | 16.2 | 22.2 | 21.8 | 23.5 | 17.3 | 22.5 | 30.2 | 23.4 |
| | Model: Qwen2.5-7B | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| **Method** | **Dataset: XLSum** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | 1% | | | | | | | | | | | 3% | | | | | | | | | | | 5% | | | | | | | | | | |
| Langs | en | es | fr | ko | pt | sw | tr | uk | vi | zh | **Avg.** | en | es | fr | ko | pt | sw | tr | uk | vi | zh | **Avg.** | en | es | fr | ko | pt | sw | tr | uk | vi | zh | **Avg.** |
| *Rand* | 25.9 | 21.9 | 28.4 | 15.5 | 23.7 | 19.9 | 20.6 | 17.9 | 20.4 | 26.6 | 22.1 | 26.2 | 21.7 | 29.0 | 15.5 | 23.0 | 18.2 | 21.9 | 17.2 | 21.1 | 28.3 | 22.2 | 26.4 | 21.4 | 29.0 | 16.0 | 23.4 | 20.2 | 21.5 | 17.6 | 21.6 | 29.6 | 22.7 |
| | *Feature-based Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *KMC\** | 27.4 | 20.7 | 29.6 | 16.9 | 24.3 | 20.3 | 21.2 | 15.8 | 19.7 | 29.1 | 22.5 | 25.9 | 22.6 | 26.5 | 15.6 | 24.9 | 18.0 | 22.5 | 17.4 | 20.3 | 26.9 | 22.0 | 26.8 | 21.4 | 26.0 | 14.1 | 23.3 | 21.3 | 19.3 | 15.8 | 20.8 | 29.4 | 21.8 |
| *MTLD\** | 25.7 | 21.8 | 27.8 | 15.0 | 21.0 | 19.7 | 21.2 | 16.7 | 21.2 | 26.7 | 21.7 | 29.6 | 20.4 | 27.7 | 17.3 | 23.0 | 19.1 | 21.9 | 16.0 | 22.0 | 29.2 | 22.6 | 28.9 | 21.7 | 28.3 | 16.1 | 23.4 | 16.2 | 21.0 | 17.9 | 21.3 | 29.4 | 22.4 |
| *Nat* | 24.1 | 22.3 | 26.7 | 14.9 | 23.3 | 21.2 | 20.3 | 13.9 | 19.0 | 28.5 | 21.4 | 24.7 | 21.5 | 25.8 | 14.3 | 22.8 | 17.4 | 19.0 | 17.0 | 18.8 | 27.2 | 20.9 | 22.9 | 21.2 | 25.5 | 15.7 | 21.3 | 18.5 | 21.2 | 15.5 | 17.8 | 25.6 | 20.5 |
| *Coh* | 24.7 | 21.5 | 26.5 | 15.7 | 24.4 | 21.2 | 21.5 | 15.6 | 21.2 | 24.7 | 21.7 | 25.9 | 21.3 | 27.3 | 14.5 | 23.7 | 21.8 | 20.9 | 15.5 | 21.3 | 26.2 | 21.8 | 24.7 | 22.4 | 27.1 | 16.1 | 24.8 | 23.3 | 21.5 | 15.8 | 19.9 | 28.8 | 22.4 |
| *Und* | 24.4 | 20.9 | 26.4 | 14.3 | 22.4 | 22.5 | 21.1 | 13.4 | 17.4 | 27.8 | 21.1 | 24.1 | 21.9 | 25.3 | 14.2 | 22.6 | 17.1 | 19.2 | 14.7 | 18.9 | 25.5 | 20.3 | 24.2 | 22.0 | 26.7 | 14.8 | 22.9 | 16.2 | 19.9 | 16.7 | 19.1 | 25.9 | 20.9 |
| | *Target-dependent Baselines* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *DSIR\** | 25.9 | 21.9 | 27.5 | 14.4 | 23.0 | 20.9 | 21.3 | 15.8 | 20.1 | 24.9 | 21.6 | 25.0 | 21.4 | 28.0 | 14.0 | 23.2 | 17.7 | 19.0 | 15.7 | 20.5 | 27.9 | 21.3 | 26.6 | 21.2 | 27.7 | 15.4 | 25.5 | 22.2 | 20.9 | 15.3 | 21.5 | 28.6 | 22.5 |
| *LESS* | 20.9 | 18.5 | 21.3 | 12.4 | 16.7 | 19.4 | 17.4 | 14.9 | 14.3 | 23.9 | 18.0 | 22.5 | 21.9 | 25.6 | 12.1 | 20.2 | 16.0 | 18.1 | 15.3 | 17.1 | 25.7 | 19.4 | 25.2 | 20.5 | 27.8 | 14.4 | 22.3 | 15.7 | 20.0 | 17.3 | 19.1 | 27.2 | 20.9 |
| | *Applying LangGPS (to random selection and three best well-performing baselines marked with \*)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| *Rand* | 26.6 | 22.1 | 28.1 | 15.1 | 24.3 | 18.6 | 20.3 | 17.9 | 19.9 | 26.5 | 22.0 | 25.3 | 21.4 | 28.5 | 15.0 | 23.7 | 18.3 | 21.4 | 17.1 | 20.4 | 26.2 | 21.7 | 25.3 | 21.3 | 29.9 | 15.9 | 23.3 | 22.5 | 21.3 | 16.9 | 20.9 | 26.3 | 22.4 |
| *KMC* | 24.4 | 21.6 | 28.1 | 14.8 | 24.1 | 16.0 | 20.1 | 17.4 | 21.2 | 28.3 | 21.6 | 29.4 | 22.6 | 28.9 | 16.5 | 23.3 | 21.8 | 21.4 | 15.1 | 22.0 | 29.1 | 23.0 | 27.8 | 21.5 | 28.6 | 16.9 | 24.3 | 23.9 | 22.4 | 15.6 | 20.6 | 28.3 | 23.0 |
| *MTLD* | 25.3 | 21.1 | 26.5 | 15.5 | 20.9 | 18.0 | 20.4 | 17.5 | 19.6 | 26.3 | 21.1 | 25.3 | 22.2 | 28.8 | 15.2 | 23.0 | 19.6 | 21.7 | 16.3 | 20.9 | 24.6 | 21.8 | 24.9 | 21.3 | 29.3 | 15.8 | 25.0 | 22.4 | 21.9 | 16.4 | 22.3 | 25.1 | 22.4 |
| *DSIR* | 24.4 | 21.5 | 25.6 | 14.1 | 23.2 | 19.4 | 21.7 | 16.5 | 19.0 | 24.3 | 21.0 | 25.1 | 22.4 | 29.0 | 15.1 | 24.2 | 20.0 | 22.8 | 16.9 | 21.5 | 24.7 | 22.2 | 25.6 | 22.1 | 28.0 | 15.7 | 23.7 | 23.0 | 23.1 | 16.4 | 21.6 | 25.6 | 22.5 |

Table 10: The detailed performance results of different language subset on *XLSum* dataset.